

TECHNICAL REPORT

CONSTRAINED ADAPTIVE DIGITAL ARRAY PROCESSING

for

National Aeronautics and Space Administration

Electronics Research Center

under

NASA Grant, NGR 33-006-040

Prepared by

Leonard P. Winkler

Mischa Schwartz

Department of Electrical Engineering

Polytechnic Institute of Brooklyn

PIBEE 71-007

May 1971



|                   |   |                                       |
|-------------------|---|---------------------------------------|
| FACILITY FORM 602 | <u>71-34162</u><br>(ACCESSION NUMBER)             | <u>                    </u><br>(THRU) |
|                   | <u>174</u><br>(PAGES)                             | <u>53</u><br>(CODE)                   |
|                   | <u>CR-121292</u><br>(NASA CR OR TMX OR AD NUMBER) | <u>07</u><br>(CATEGORY)               |
|                   |   |                                       |

TECHNICAL REPORT

CONSTRAINED ADAPTIVE DIGITAL ARRAY PROCESSING

for

National Aeronautics and Space Administration

Electronics Research Center

under

NASA Grant NGR 33-006-040

Prepared by

Leonard P. Winkler

Mischa Schwartz

Department of Electrical Engineering

Polytechnic Institute of Brooklyn

PIBEE 71-007

May 1971

## ABSTRACT

This investigation is concerned with automatically making an array of detectors form a beam in a desired direction in space when unknown interfering noise is present so as to maximize the output signal-to-noise ratio (SNR) subject to a constraint on the super-gain ratio (Q-factor). Tapped delay line structures combined with iterative gradient techniques to adjust the tap weights are used to do this. }

First, we investigate the relationship between viewing the detectors as a "detector array" and viewing the detectors as a "multichannel filter."

Next, starting from the multichannel filter point of view we investigate the sensitivity of the SNR to random errors in the tap weight settings and random errors in our knowledge of the detector locations. Because this calculation is exceedingly difficult from the multichannel filter approach, we will use the previously derived relationship to show that this sensitivity is essentially given by the super-gain ratio. We show that when we use linear arrays of detectors separated by one-half wavelength or less, this sensitivity factor may become very large when we use those currents and phases (or tap weights) which maximize the SNR, thus indicating that we should not try to design our detector pattern or multichannel filter coefficients on the basis of maximizing the SNR alone, but rather on the basis of maximizing the SNR subject to a constraint on the super-gain ratio.

We then develop a computationally fast numerical method of finding the optimum excitations which maximize the SNR subject to a super-gain ratio constraint when the interfering noise is known.

Next, we try to analytically consider adaptive algorithms which maximize the SNR subject to a constraint on the super-gain ratio when unknown interfering noise is present, but because the SNR and super-gain ratio are nonlinear quantities, it turns out to be exceedingly difficult to prove convergence of the algorithms to the optimal solution, or to find the algorithms' rates of convergence. Thus, solely for the purpose of mathematical tractability, we consider adaptive algorithms which minimize the mean square error (MSE) subject to a linear constraint.

Finally we present the results of computer simulations of algorithms which maximize the SNR subject to a constraint on the super-gain ratio when unknown interfering noise is present.

## TABLE OF CONTENTS

|   | <u>Page</u> |
|---|-------------|
| <u>Chapter 1: Introduction</u>  | 1           |
| <u>Chapter 2: Equivalence Between "Detector Pattern" and "Multichannel Filter" Viewpoints in Designing Optimum Arrays</u>                         | 7           |
| Section 2.1: "Detector Pattern" Approach  | 8           |
| Section 2.2: "Multichannel Filter" Approach   | 12          |
| Section 2.3: Relationships Between the "Detector Pattern" and "Multichannel Filter" Approaches  | 17          |
| Appendix A: Maximization of the SNR   | 27          |
| Appendix B: Evaluation of $\phi_n(\tau, \underline{x}_k - \underline{x}_l)$ for Temporally Monochromatic and White Noise                          | 29          |
| Appendix C: Evaluation of the A Matrix  | 31          |
| Appendix D: Evaluation of the Q Matrix  | 33          |
| <u>Chapter 3: Error Analysis of Point Detector Arrays</u>   | 34          |
| Section 3.1: Sensitivity of the SNR to Random Errors in the Detector Excitations and Locations  | 35          |
| Section 3.2: Maximization of the SNR Subject to a Constraint on the Super-Gain Ratio  | 53          |
| Appendix A: Statistical Formulation of the Super-Gain Ratio   | 60          |
| Appendix B: Maximization of the SNR Subject to a Constraint   | 67          |
| <u>Chapter 4: Minimization of the MSE Subject to One Linear Constraint</u>  | 73          |
| Section 4.1: Derivation of MSE and Linear Constraint Equation   | 74          |
| Section 4.2: The Analytic (Lagrange) Solution   | 77          |
| Section 4.3: Use of the Projected Gradient Algorithm to Adaptively Adjust the Tap Weights   | 80          |
| Section 4.3.1: The Algorithm, Proof of Convergence, and Bounds on the Rate of Convergence, if the Gradient is Known                               | 81          |
| Section 4.3.2: The Algorithm, Proof of Convergence, and Bounds on the Rate of Convergence if the Gradient is Estimated                            | 87          |
| Section 4.3.3: The Algorithm, Proof of Convergence and Bounds on the Rate of Convergence if the Gradient is Estimated, and the Estimate is Noisy. | 92          |
| Section 4.4: Computer Simulations   | 96          |

## Table of Contents-continued

|  | <u>Page</u> |
|--|-------------|
| Appendix A: Proof of Convergence and Bounds on the Asymptotic Variance   | 115         |
| Appendix B: Rosen's Gradient Projection Algorithm  | 120         |
| <u>Chapter 5: Adaptive Algorithm to Minimize MSE Subject to a "Soft" Constraint</u>  | <u>126</u>  |
| Section 5.1: Introduction  | 126         |
| Section 5.2.1: The Algorithm, Proof of Convergence, and Bounds on the Rate of Convergence if the Gradient is Known                               | 129         |
| Section 5.2.2: The Algorithm, Proof of Convergence, and Bounds on the Rate of Convergence if the Gradient is Estimated                           | 133         |
| Section 5.2.3: The Algorithm, Proof of Convergence, and Bounds on the Rate of Convergence if the Gradient is Estimate, and the Estimate is Noisy | 137         |
| <u>Chapter 6: Computer Simulations of Nonlinear Problem and Conclusions</u>  | <u>143</u>  |
| Section 6.1: Antenna Theory Approach   | 144         |
| Section 6.2: Multichannel Filter Approach  | 148         |
| Section 6.3: Maximization of SNR Subject to $Q \leq q$   | 153         |
| Section 6.4: The Gradient Projection Algorithm   | 155         |
| Section 6.5: Conclusions   | 165         |

# TABLE OF FIGURES AND GRAPHS

|  | <u>Page</u> |
|--|-------------|
| 1.1      Convergence of an Arbitrary Tap Weight to its<br>Steady-State Value | 2           |
| 2.1.1    Detector Array  | 8           |
| 2.2.1    Multichannel Filter Structure                                       | 12          |
| 2.2.2    Incident Signal Field   | 14          |
| 2.3.1    Incident Noise Field  | 18          |
| 2.3.2    Correlation between two Detectors                                   | 21          |
| 2.3.3    Detector Array  | 25          |
| 3.1.1    Typical Power Pattern   | 36          |
| 3.1.2    Four Element Linear Array   | 38          |
| 3.1.3    Ten Element Linear Array  | 42          |
| 3.1.4    Four Element Array - Broadside Signal                               | 43          |
| 3.1.5    Four Element Array - Broadside Signal                               | 44          |
| 3.1.6    Ten Element Array - Broadside Signal                                | 45          |
| 3.1.7    Ten Element Array - Broadside Signal                                | 46          |
| 3.1.8    Four Element Array - Endfire Signal                                 | 47          |
| 3.1.9    Four Element Array - Endfire Signal                                 | 48          |
| 3.1.10   Ten Element Array - Endfire Signal                                  | 49          |
| 3.1.11   Ten Element Array - Endfire Signal                                  | 50          |
| 3.1.12   Extension of Fig. 3.1.4   | 51          |
| 4.1.1    Processor Configuration   | 74          |
| 4.2.1    Typical MSE Level Curves and Constraint                             | 78          |
| 4.3.1    Intuitive Idea behind Projected Gradient Algorithm                  | 80          |
| 4.3.2 $\xi$ vs. $k$  | 86          |
| 4.3.3    Bounds on $k_{\max}$  | 87          |
| 4.4.1    Gradient Known, No Additive Noise                                   | 100         |
| 4.4.2    Gradient Known, No Additive Noise                                   | 101         |

Table of Figures and Graphs Continued:

|  | <u>Page</u> |
|--|-------------|
| 4.4.3 Gradient Estimated, No Additive Noise                      | 103         |
| 4.4.4 Gradient Estimated, No Additive Noise                      | 104         |
| 4.4.5 Gradient Estimated, Plus Additive Noise                    | 107         |
| 4.4.6 Gradient Estimated, Plus Additive Noise                    | 108         |
| 4.4.7 Gradient Estimated, Plus Additive Noise                    | 109         |
| 4.4.8 Gradient Estimated, Plus Additive Noise                    | 110         |
| 4.4.9 Gradient Estimated, Plus Additive Noise                    | 111         |
| 4.4.10 Gradient Estimated, Plus Additive Noise                   | 112         |
| 4.4.11 Gradient Estimated, Plus Additive Noise                   | 113         |
| 4.4.12 Gradient Estimated, Plus Additive Noise                   | 114         |
| B1 Diagram for Example One                                       | 124         |
| B2 Diagram for Example Two                                       | 125         |
| 5.1.1 Constraint and Penalty Function Level Curves               | 126         |
| 6.2.1 Processor Structure  | 149         |
| 6.4.1 Gradient Projection Operation                              | 156         |
| 6.4.2 Broadside Gradient Known, No Additive Detector Noise       | 159         |
| 6.4.3 Broadside Gradient Estimated, No Additive Detector Noise   | 160         |
| 6.4.4 Broadside Gradient Estimated, Plus Additive Detector Noise | 161         |
| 6.4.5 Endfire Gradient Known, No Additive Detector Noise         | 162         |
| 6.4.6 Endfire Gradient Estimated, No Additive Detector Noise     | 163         |
| 6.4.7 Endfire Gradient Estimated, Plus Additive Detector Noise   | 164         |



## CHAPTER 1

### INTRODUCTION

This investigation is concerned with the optimal design of a detector array and signal processor to maximize the output signal-to-noise ratio (SNR) subject to a constraint on the super-gain ratio (Q-factor). We will present and analyze an iterative gradient projection technique to achieve this optimal design even when the noise statistics are unknown to the designer a priori.

Some of the motivations for undertaking our study at the present time are:

1. The recent ability to approximate the sophisticated processing required through the use of fast, special-purpose digital computers.
2. The recent use of channels, such as are present in spacecraft and underwater communications, where the additive noise from spatially distributed noise sources predominates over the additive receiver noise.
3. The recent use of acoustic and seismic channels where the low signal frequencies used result in long signal and noise wavelengths (relative to array size), thus to high correlations between the noise at the array elements, which in turn implies that we might achieve improved performance through the use of array processing techniques.
4. The limited ability of design procedures based upon the classical concept of an antenna pattern to adequately satisfy the criteria of minimum probability of error or minimum mean squared error or maximum SNR.

The first three factors are self-explanatory. The last one deserves some comment. Some of the advantages (and limitations) of the classical antenna pattern approach to the design of array processors are:

1. The approach subdivides the system design problem into two separate pieces. An antenna engineer designs the array (spatial processor) and independently, a communications engineer takes the single channel antenna output and designs the temporal processor to give, for example, the best, in some sense, estimate of the transmitted signal.

This would seem to be an advantage, however, Gaarder<sup>(2)-(3)</sup> has shown that this factoring of the optimum processor into spatial and temporal processors is, in general, impossible, and consequently, processors designed on this principle are suboptimum.

2. The concept of an antenna pattern assumes that we are dealing with monochromatic or quasi-monochromatic fields. For the wideband signals coming into use, there is no easy way of combining the various frequency components together.

Previous researchers<sup>(1)-(11)</sup> have considered the design of detector arrays to maximize some criterion without constraints, both from the "detector pattern" point of view and from the "multichannel filter" point of view. More recently<sup>(12)-(18)</sup> investigators have devised adaptive algorithms to enable processing structure composed of tapped delay lines (such as that shown in Fig. 6.2.1) to converge to an optimal structure even when the noise statistics are unknown to the designer a priori. These algorithms are similar to those used to adaptively equalize telephone and other dispersive communication channels.

These previous authors have designed adaptive algorithms which minimized the MSE, or maximized the SNR, by using iterative gradient techniques to make the tap weights converge to values which optimize the MSE or SNR in the steady state. Any individual tap weight usually converges to its steady-state value in a manner similar to that shown in Fig. 1.1 below.

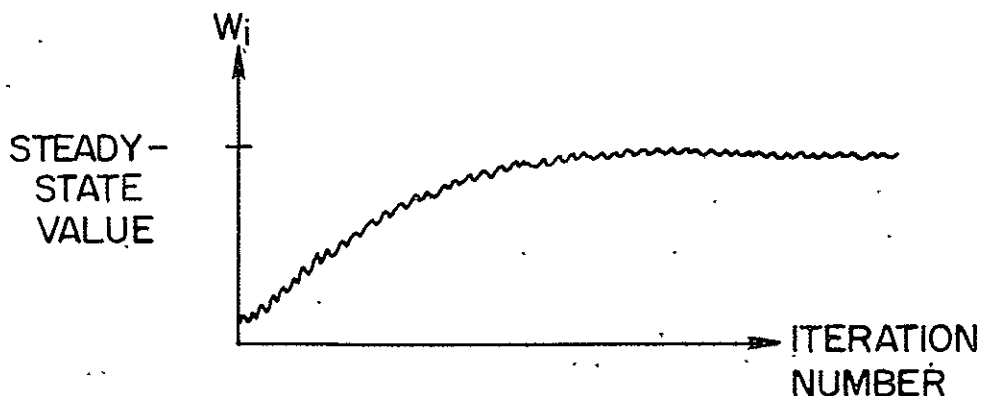


Fig. 1.1 Convergence of an arbitrary tap weight to its steady-state value

In the steady state, each tap weight can be viewed as having a nominal value plus a random variation about this nominal value. If we use the unbiased algorithms of Widrow, (12, 13) Griffiths (15) and Somin (14) the nominal value is the same as the optimal value of the tap weight. However a question that immediately arises is the following: How sensitive is the SNR to the small random variations in the tap weights about their nominal values?

In chapter three we will show that, depending upon the geometry of the detector array, the SNR can be very sensitive to these small random variations, and we will derive an expression for this sensitivity.

In order to derive the expression for the sensitivity, some reformulation of what previous investigators have done, both from the "detector pattern" point of view and from the "multichannel filter" point of view, will be necessary. This will be covered in chapter two where we will also demonstrate that both approaches lead to the same results under a monochromatic assumption, which is to be expected, since there is only one physical problem. The reason for our reformulation is as follows: We will be able to express the SNR in the form  $\frac{\underline{Z}^* \underline{P} \underline{Z}}{\underline{Z}^* \underline{Q} \underline{Z}}$  or  $\frac{\underline{I}^* \underline{C} \underline{I}}{\underline{I}^* \underline{A} \underline{I}}$  where the vector  $\underline{Z}$  represents the complex gains (or tap weights) in the multichannel filter approach and the vector  $\underline{I}$  represents the excitation currents in the detector pattern. By the sensitivity of the SNR to random errors in the tap weights we mean that if we replace  $\underline{Z}$  by  $\underline{Z}_N + \underline{Z}_R$  where  $N$  denotes the nominal value and  $R$  denotes the random fluctuations about this nominal value, the expected value of  $\frac{\underline{Z}^* \underline{P} \underline{Z}}{\underline{Z}^* \underline{Q} \underline{Z}}$

may turn out to be of the form  $E \left\{ \frac{\underline{Z}^* \underline{P} \underline{Z}}{\underline{Z}^* \underline{Q} \underline{Z}} \right\} = \frac{\underline{Z}_N^* \underline{P} \underline{Z}_N}{\underline{Z}_N^* \underline{Q} \underline{Z}_N} + \text{an additional term,}$  and we then define the ratio of the additional term to the nominal term as our sensitivity factor. However, using this approach, the calculation of  $E \left\{ \frac{\underline{Z}^* \underline{P} \underline{Z}}{\underline{Z}^* \underline{Q} \underline{Z}} \right\}$  is exceedingly complex. Instead, because we showed in chapter two that the detector pattern and multichannel filter approaches were interchangeable, we will use the detector pattern approach and rewrite the SNR expression above in terms of the power pattern, which in turn depends upon the excitation currents, and then by examining a picture of a typical power pattern, we will be lead by physical reasoning to approximate the sensitivity of the SNR to random variations in the tap weights, by the super-gain ratio,

which is a measure of the sensitivity of the power at the peak of the beam to random errors in the detector excitations. In other words, instead of saying that changes in the tap weights cause changes in the SNR, we are now saying that changes in the tap weights cause changes in the peak of the power pattern which in turn is the main reason the SNR changes. Thus if we constrain changes in the peak of the power pattern we will also automatically constrain changes in the SNR. The advantage is that we can easily derive an expression for changes in the peak of the power pattern due to changes in the tap weights (or detector currents), whereas we cannot easily derive an expression for changes in the SNR due to changes in the tap weights.

As mentioned before, we will show in chapter three that although, for a particular array geometry (specifically a linear array of detectors separated by half a wavelength, where the signal is impinging from endfire), we might initially be lead to believe that we can achieve very good performance by setting (usually by means of an adaptive algorithm) the tap weights equal to those values which maximize the SNR, if we also look at the super-gain ratio, we will see that in practice we will not get this good performance because of the extreme sensitivity of the SNR to the small deviations in the tap weights from their optimal values.

After demonstrating this, section 3.2 goes on to answer the question of how high a SNR can we get if we constrain the super-gain ratio to equal some reasonable value. In order to do this we will extend the work of Lo, Lee and Lee, <sup>(19)</sup> who recently developed a numerical method of solving this problem. Our contribution makes use of a state variable technique which enables us to reduce the numerical problem from one of finding the complex roots of a high order polynomial with complex coefficients (in all the specific numerical cases treated in the paper by Lo, Lee and Lee the coefficients of the polynomials were real, but this is not necessarily true in general) to one of finding eigenvalues of a real matrix which is considerably faster and easier to do.

Next, we tried to analytically consider adaptive algorithms which would maximize the SNR subject to a constraint on the super-gain ratio when unknown interfering noise is present. Because the SNR and super-gain ratio are nonlinear quantities, it turned out to be exceedingly difficult to prove convergence of the algorithms to the optimum solution, or to find the algorithms' rates of convergence. Thus, solely for the purpose of mathematical tractability (the actual nonlinear problem will be simulated on a computer in

chapter six to obtain some numerical indication of convergence and convergence rates), chapter four analyzes an adaptive projection algorithm which minimizes the mean square error (MSE) subject to a linear constraint. We prove that an algorithm of the form

$$\underline{W}_{j+1} = \underline{W}_j - k P \nabla_{\underline{W}_j} (\text{MSE})$$

converges to the Lagrange solution in real-time, with an easily expressible bound on the convergence rate. Here  $k$  is the step size,  $P$  is a matrix projection operator (20)-(21) and  $\nabla_{\underline{W}_j}$  is the gradient of the MSE with respect to  $\underline{W}_j$ . We also proved convergence and found bounds on the rate of convergence when  $\nabla_{\underline{W}_j} (\text{MSE})$  was (1) known exactly (2) estimated, and (3) estimated by a noisy estimate. Physically these cases correspond to (1) knowing the interfering noise field exactly (2) using the instantaneous values of the noise that are present at the outputs of the detectors (or at the outputs of each of the delay elements comprising our tapped delay lines) as estimates of the noise correlation matrix, e. g. replacing  $E \{n_i(t)n_j(t)\}$  by  $n_i(t_k)n_j(t_k)$  at iteration  $k$ , and (3) accounting for self-noise in the detectors and tapped delay lines by replacing  $E \{n_i(t)n_j(t)\}$  by  $n_i(t_k)n_j(t_k) + \xi_k$  at iteration  $k$  where  $\xi_k$  is additive white gaussian noise.

Chapter five is an investigation of an adaptive penalty algorithm to minimize the MSE subject to a linear constraint. Specifically we prove that algorithms of the form

$$\underline{W}_{j+1} = \underline{W}_j - k \nabla_{\underline{W}_j} \left( \text{MSE} + K_1 \left[ \underline{W}_j^T \underline{n}_1 - a \right]^2 \right)$$

where  $\underline{W}^T \underline{n}_1 - a$  is the equation defining the linear constraint, converge to the Lagrange solution of chapter four if  $K_1$  is infinite. For  $K_1$  finite, a bias is found to exist, and is investigated, along with bounds on the rates of convergence of these algorithms to their steady-state values. Again we considered the same three ways of evaluating  $\nabla_{\underline{W}_j} (\text{MSE})$ .

In chapter six, we set up and present the results of a computer simulation of the gradient projection algorithm which adaptively maximizes the SNR subject to a constraint on the super-gain ratio. We then conclude that when designing adaptive array processors one should either

1. Calculate the super-gain ratio for the geometry under consideration for all possible incident signal directions and if we are sure that the

super-gain ratio can never become intolerably high feel free to use the adaptive gradient algorithms proposed by previous authors, or

2. Use the constrained adaptive algorithms developed in this investigation, which will assure us that we get the highest SNR possible subject to a constraint on the super-gain ratio should the value of the super-gain ratio exceed some preset value we have chosen.

## CHAPTER 2

### Equivalence Between "Detector Pattern" and "Multichannel Filter" Viewpoints in Designing Optimum Arrays

In this chapter, we will consider the following problem: Given an array of point detectors at known locations in space, how should we "design" the array so as to maximize the output SNR? This problem has been solved before—as a matter of fact, it has been solved twice before, once by antenna engineers, who solved for those detector current excitations which maximized the SNR through the use of the "detector pattern" concept, and again by communication engineers who viewed the array as a multichannel filter and solved for those filter coefficients which maximized the SNR, through the use of statistical quantities such as the covariances of the signal and noise fields.

As explained in more detail in the first chapter, we will reformulate what these previous investigations have done, and show that the two approaches are equivalent (i.e., lead to the same optimum value of the SNR under a monochromatic noise assumption) in order that we may, in chapter three, easily switch from the multichannel filter point of view to the detector pattern viewpoint when evaluating the sensitivity of the SNR to small random variations in the tap weights.

In section 2.1 we derive the optimum currents and the resulting value of the SNR when these currents are used to excite the detector array. All our results will be a function of the assumed incident noise power. In section 2.2 we derive the optimum filter coefficients and the resulting value of the SNR when these filter coefficients are used in the multichannel filter. These results will be a function of the assumed noise space-time correlation function. In section 2.3 we will express the space-time correlation function used in section 2.2 as a direct function of the incident noise power used in section 2.1 and then show that under the monochromatic noise assumption, the detector pattern approach and the multichannel filter approach, yield exactly the same value of the SNR, and moreover, we will be able to see that the currents of section 2.1 correspond to the filter coefficients of section 2.2. This analogy will be used in the following chapter to construct a quantity which is defined in terms of communication theory quantities (e.g., covariance), and corresponds to the super-gain ratio of antenna theory.

## Section 2.1 "Detector Pattern" Approach

The material in this section follows the approach of Lo, Lee and Lee.<sup>(19)</sup>

Assume we have  $N$  isotropic detectors located at arbitrary positions in space, specified by Cartesian coordinates  $\underline{r}_n = (x_n, y_n, z_n)$  as shown in Fig. 2.1.1.

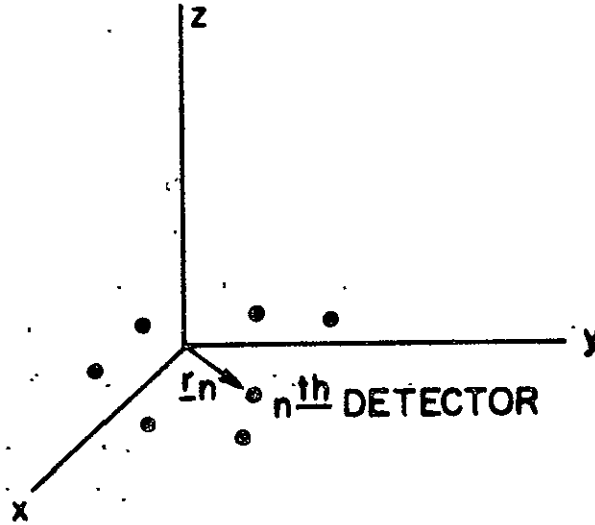


Fig. 2.1.1 Detector Array

The current in the  $n^{\text{th}}$  detector will be denoted by  $I_n$ . Let us define

$$\underline{I}^* \equiv (I_1, I_2, \dots, I_N) \quad (2.1.1)$$

where the asterisk denotes adjoint. The detector pattern is given by

$$p(\theta, \phi) = \sum_{n=1}^N I_n e^{jk \underline{r}_0 \cdot \underline{r}_n} \quad (2.1.2)$$

where the  $\underline{r}$ 's are given by

$$\underline{r}_0 = \sin \theta \cos \phi \underline{x}_0 + \sin \theta \sin \phi \underline{y}_0 + \cos \theta \underline{z}_0$$

$$\underline{r}_n = x_n \underline{x}_0 + y_n \underline{y}_0 + z_n \underline{z}_0 = \text{the position of the } n^{\text{th}} \text{ element}$$

Since  $k = \frac{2\pi}{\lambda}$  we have

$$k \underline{r}_0 \cdot \underline{r}_n = 2\pi \left[ \frac{x_n}{\lambda} \sin \theta \cos \phi + \frac{y_n}{\lambda} \sin \theta \sin \phi + \frac{z_n}{\lambda} \cos \theta \right]$$



We will define

$$\psi_n \equiv k \underline{r}_o \cdot \underline{r}_n = 2\pi \left[ \frac{x_n}{\lambda} \sin \theta \cos \phi + \frac{y_n}{\lambda} \sin \theta \sin \phi + \frac{z_n}{\lambda} \cos \theta \right] \quad (2.1.3)$$

Equation (2.1.2) becomes

$$p(\theta, \phi) = \sum_{n=1}^N I_n e^{j\psi_n} \equiv \underline{I}^* \underline{V} \quad (2.1.4)$$

where  $\underline{V}$  is given by

$$\underline{V} = \begin{bmatrix} e^{+j\psi_1} \\ \vdots \\ e^{+j\psi_n} \end{bmatrix} \quad (2.1.5)$$

If we assume the normalized signal is incident from direction  $(\theta_o, \phi_o)$ , then the received signal power is given by

$$\begin{aligned} S &= \iint_{\theta, \phi} |p(\theta, \phi)|^2 \delta(\theta - \theta_o, \phi - \phi_o) d\Omega \\ &= [\underline{I}^* \underline{V}_1]^2 = \underline{I}^* \underline{V}_1 \underline{V}_1^* \underline{I} \end{aligned} \quad (2.1.6)$$

where

$$\underline{V}_1^* = [e^{-j\psi_1^o} \dots e^{-j\psi_n^o}] \quad (2.1.7)$$

$$\text{and } \psi_n^o = 2\pi \left[ \frac{x_n}{\lambda} \sin \theta_o \cos \phi_o + \frac{y_n}{\lambda} \sin \theta_o \sin \phi_o + \frac{z_n}{\lambda} \cos \theta_o \right] \quad (2.1.8)$$

Define the matrix C by

$$\underline{V}_1 \underline{V}_1^* = \begin{bmatrix} e^{j\psi_1^o} \\ \vdots \\ e^{j\psi_n^o} \end{bmatrix} \begin{bmatrix} e^{-j\psi_1^o} & \dots & e^{-j\psi_n^o} \end{bmatrix} = C \quad (2.1.9)$$

Note that C is a Hermitian positive definite matrix (dyadic)

$$\text{Proof: } \underline{x}^* C \underline{x} = \underline{x}^* \underline{V}_1 \underline{V}_1^* \underline{x} = |\underline{x}^* \underline{V}_1|^2 > 0 \text{ if } \underline{x} \neq 0$$

Thus

$$S = \underline{I}^* C \underline{I} \quad (2.1.10)$$

Let us assume that the spatial distribution of the noise power is given by  $T(\theta, \phi)$ . Then the noise power received is:

$$\begin{aligned} N &= \int_{\theta} \int_{\phi} |P(\theta, \phi)|^2 T(\theta, \phi) d\Omega \\ &= \int_{\theta} \int_{\phi} \underline{I}^* \underline{V} \underline{V}^* \underline{I} T(\theta, \phi) d\Omega \end{aligned} \quad (2.1.11)$$

Since the currents  $\underline{I}_n$  are not functions of  $\theta$  or  $\phi$

$$N = \underline{I}^* \left[ \int_{\theta} \int_{\phi} \underline{V} \underline{V}^* T(\theta, \phi) d\Omega \right] \underline{I}$$

Define the matrix A by

$$N = \underline{I}^* \underline{A} \underline{I} \quad (2.1.12)$$

where the elements of the matrix A are given by  $a_{ij}$

$$a_{kl} = \int_{\theta} \int_{\phi} e^{+j\psi_k} e^{-j\psi_l} T(\theta, \phi) d\Omega$$

The matrix A is positive definite

$$\text{Proof: } \underline{x}^* \underline{A} \underline{x} = \underline{x}^* \int_{\theta} \int_{\phi} \underline{V} \underline{V}^* T(\theta, \phi) d\Omega \underline{x}$$

$$= \int_{\theta} \int_{\phi} \left[ \underline{x}^* \underline{V} \right] \left[ \underline{V}^* \underline{x} \right] T(\theta, \phi) d\Omega$$

Because  $T(\theta, \phi)$  is always positive, we may write it as

$$T(\theta, \phi) = g(\theta, \phi) g^*(\theta, \phi) \text{ where } g \text{ and } g^* \text{ are scalars}$$

Thus

$$\underline{x}^* \underline{A} \underline{x} = \int_{\theta} \int_{\phi} \left[ g \underline{x}^* \underline{V} \right] \left[ g^* \underline{V}^* \underline{x} \right] d\Omega$$

$$= \int \int_{\theta \phi} |g \underline{x}^* \underline{V}|^2 d\Omega$$

Since the integrand is positive

$$\underline{x}^* A \underline{x} > 0 \text{ if } \underline{x} \neq 0$$

Q E D

The signal-to-noise ratio (SNR) is then

$$\text{SNR} = \frac{\underline{I}^* C \underline{I}}{\underline{I}^* A \underline{I}} \quad (2.1.13)$$

We may use the calculus of variations to find the value of  $\underline{I}$  which maximizes the SNR. From Appendix A

$$\underline{I}_{\text{optimum}} = A^{-1} \underline{V}_1 \quad (2.1.14)$$

The value of the SNR when  $\underline{I} = \underline{I}_{\text{opt}}$  is

$$\text{SNR} = \frac{\underline{I}_{\text{opt}}^* C \underline{I}_{\text{opt}}}{\underline{I}_{\text{opt}}^* A \underline{I}_{\text{opt}}} = \underline{V}_1^* A^{-1} \underline{V}_1$$

The best SNR that we can achieve by using the "detector pattern" approach to the problem of optimizing the SNR is thus

$$\text{SNR} = \underline{V}_1^* A^{-1} \underline{V}_1 \quad (2.1.15)$$

We will now find an expression for the best SNR we can achieve by using the multichannel filter approach to the problem of optimizing the SNR and then show under what conditions the two approaches yield the same value for the best SNR.

## Section 2.2 Multichannel Filter Approach

Assuming that we know the noise space-time correlation function, let us now find the optimum multichannel filter, optimum in the sense that we will find the  $z_i$ 's (see Fig. 2.2.1) which maximize the SNR. Once the coefficients of the optimum filter have been found, we will be able to write an expression for the best SNR we can achieve through the use of the multichannel filter approach.

The material in this section follows the approach of Edelblute, Fisk and Kinnison<sup>(8)</sup>.

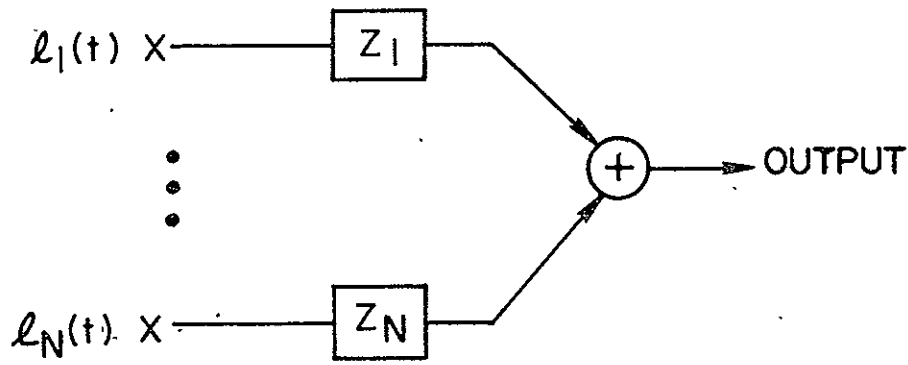


Fig. 2.2.1 Multichannel filter structure

The SNR at the multichannel filter output when  $l_i(t) = s_i(t) + n_i(t)$  is received is given (under the assumption that the signal and noise are complex uncorrelated random waveforms) by

$$\text{SNR} = \frac{\sum_{i=1}^N \sum_{j=1}^N z_i^* z_j p_{ij}}{\sum_{i=1}^N \sum_{j=1}^N z_i^* z_j q_{ij}} = \frac{\underline{Z}^* \underline{P} \underline{Z}}{\underline{Z}^* \underline{Q} \underline{Z}} \quad (2.2.1)$$

$$\text{where } E \{ s_i^*(t) n_j(t) \} = E \{ n_i^*(t) s_j(t) \} = 0 \quad \forall_{ij} \quad (2.2.2)$$

$$p_{ij} \equiv E \{ s_i^*(t) s_j(t) \} \quad (2.2.3)$$

$$q_{ij} \equiv E \{ n_i^*(t) n_j(t) \} \quad (2.2.4)$$

$$\underline{Z} = \begin{bmatrix} z_1 \\ \vdots \\ z_N \end{bmatrix} \quad (2.2.5)$$

Note that P and Q are correlation matrices and thus are Hermitian positive semidefinite (we will assume that Q is positive definite, which is generally true in practice - the Q matrix is usually of the form  $Q = \alpha I + \tilde{Q}$  where the  $\alpha I$  term is due to additive self-noise at each detector, thus guaranteeing the existence of  $Q^{-1}$ )

Note the similarity between equation (2.2.1) and equation (2.1.13). Also note that the SNR is independent of the magnitude of  $\underline{Z}$ . Let us now find the value of  $\underline{Z}$  that maximizes the SNR by using the calculus of variations, i. e.

$$\text{maximize } L = \frac{\underline{Z}^* P \underline{Z}}{\underline{Z}^* Q \underline{Z}} \quad (2.2.6)$$

This equation is of the same form as equation (A1) of Appendix A.

By the same reasoning as in section 2.1 (see equation 2.1.15) we have

$$P \underline{Z}_o \left( \underline{Z}_o^* Q \underline{Z}_o \right) - Q \underline{Z}_o \left( \underline{Z}_o^* P \underline{Z}_o \right) = 0 \quad (2.2.7)$$

where  $\underline{Z}_o$  = optimum  $\underline{Z}$

$$P \underline{Z}_o = \underbrace{\left[ \frac{(\underline{Z}_o^* P \underline{Z}_o)}{(\underline{Z}_o^* Q \underline{Z}_o)} \right]}_{\text{scalar}} Q \underline{Z}_o$$

$$\text{Let } G_o \equiv \frac{(\underline{Z}_o^* P \underline{Z}_o)}{(\underline{Z}_o^* Q \underline{Z}_o)} \quad (2.2.8)$$

Thus

$$P \underline{Z}_o = G_o Q \underline{Z}_o \quad (2.2.9)$$

Equation (2.2.9) is an equation which  $\underline{Z}_0$  must satisfy, it is not however, an explicit expression for  $\underline{Z}_0$ . Motivated by this need, and seeing from section 2.1 that one way to find such an explicit expression for  $\underline{Z}_0$  is by letting the P matrix be written as  $P = \underline{U}_1 \underline{U}_1^*$  (i.e. let P be of rank 1) let us do the following:

Assume the signal field is produced by a single source located at  $(\theta_0, \phi_0)$  in the far field, which is generating a statistically known random output.

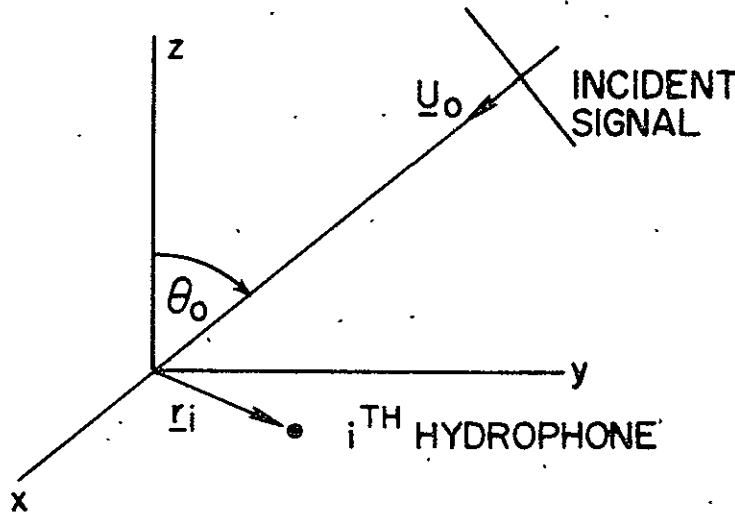


Fig. 2.2.2 Incident Signal Field

The signal may be represented in the form (where we have suppressed the  $e^{+j\omega t}$  time dependence)

$$s(\underline{x}, t) = \underline{s}(t) e^{-j \underline{k} \cdot \underline{x}} \quad \text{where } \underline{k} = \underline{u}_0 \frac{\omega}{c} = \underline{u}_0 \frac{2\pi}{\lambda}$$

At the various hydrophone locations, the received signal is

$$\begin{aligned} s(\underline{r}_i, t) &= \underline{s}(t) e^{-j \frac{\omega}{c} \underline{u}_0 \cdot \underline{r}_i} \\ &= \underline{s}(t) e^{j\omega \left( - \frac{\underline{u}_0 \cdot \underline{r}_i}{c} \right)} \end{aligned}$$

$$\text{let } \tau_i = \frac{\underline{u}_0 \cdot \underline{r}_i}{c}$$

Thus

$$s(\underline{r}_i, t) = \underline{s}(t) e^{-j\omega \tau_i} \quad (2.2.11)$$

The average signal power present in any hydrophone due to this signal is

$$\begin{aligned} S &= E \left\{ s^*(\underline{r}_i, t) s(\underline{r}_i, t) \right\} \\ &= E \left\{ \underline{s}^*(t) e^{j\omega \tau_i} \underline{s}(t) e^{-j\omega \tau_i} \right\} \\ &= E \left\{ \underline{s}^*(t) \underline{s}(t) \right\} = R_s(0) \end{aligned}$$

The normalized signal correlation matrix elements are

$$\begin{aligned} P_{ij} &= \frac{1}{R_s(0)} E \left\{ s^*(\underline{r}_i, t) s(\underline{r}_j, t) \right\} \\ &= \frac{1}{R_s(0)} E \left\{ \underline{s}^*(t) e^{j\omega \tau_i} \underline{s}(t) e^{-j\omega \tau_j} \right\} \\ &= \frac{1}{R_s(0)} e^{j\omega (\tau_i - \tau_j)} E \left\{ \underline{s}^*(t) \underline{s}(t) \right\} \\ &= e^{j\omega (\tau_i - \tau_j)} \quad (2.2.12) \end{aligned}$$

Define

$$\underline{U}_1 = \begin{bmatrix} e^{+j\omega \tau_1} \\ \vdots \\ e^{+j\omega \tau_n} \end{bmatrix} \quad \underline{U}_1^* = \begin{bmatrix} e^{-j\omega \tau_1} & \dots & e^{-j\omega \tau_n} \end{bmatrix} \quad (2.2.13)$$

Thus

$$P = \underline{U}_1 \underline{U}_1^* \quad (2.2.14)$$

We can repeat the steps leading to equation (A 3) of Appendix A, to get

$$\underline{Z}_o = \frac{(\underline{Z}_o^* Q \underline{Z}_o)}{(\underline{Z}_o^* \underline{U}_1)} Q^{-1} \underline{U}_1 \quad (2.2.15)$$

Since the SNR is independent of the magnitude of  $\underline{Z}_o$ , we see that

$$\underline{Z}_o = Q^{-1} \underline{U}_1 \quad (2.2.16)$$

is the solution for the optimum  $\underline{Z}$ .

Using this value of  $\underline{Z}$ , the optimum value of the SNR is

$$\text{SNR} = \underline{U}_1^* Q^{-1} \underline{U}_1 \quad (2.2.17)$$

This expression represents the best SNR that we can achieve by using the multichannel filter approach to the problem of optimizing the SNR.

In the next section we will investigate under what conditions this expression and the expression derived in section 2.1 for the best SNR we can achieve by using the detector pattern approach yield the same values for the optimum SNR.



### Section 2.3 Relationships between the "detector pattern" and multichannel filter approaches

In section 2.1 we found an expression for the best SNR we can achieve by using the "detector pattern" approach. In section 2.2 we found an expression for the best SNR we can achieve by using the multichannel filter approach.

We will now show that these two expressions for the optimum SNR are equivalent if the noise is monochromatic. The monochromatic assumption must be added to the multichannel filter approach because it is already inherently contained in the detector pattern approach, i.e. in deriving equation (2.1.2) the detector excitations were assumed to be monochromatic.

Showing that the two SNR expressions are equivalent entails expressing the space-time correlation functions  $\phi_n(\tau, \underline{x}_k, \underline{x}_l)$  used in section 2.2 (i.e. used in the sense that  $q_{kl} = E \left\{ n_k^*(t) n_l(t) \right\} = \phi_n(0, \underline{x}_k - \underline{x}_l)$ ) as direct functions of the incident noise power  $T(\theta, \phi)$  used in section 2.1. In order to do this we will first find the space-time correlation functions between the point detectors in the array as functions of the incident noise field. Next we will find the incident noise power as a function of the incident noise field. Finally we will be able to express the space-time correlation functions as direct functions of the incident noise power.

We will then apply the general theory to certain special noise power distributions and a particular array configuration. We will show, that under a monochromatic noise assumption, for these noise power distributions and this array configuration, the detector pattern approach and the multichannel filter approach yield exactly the same SNR results. Although we have used particular noise power distributions and a particular array configuration, this was only done to simplify the evaluation of certain integrals, and the equivalence does not depend upon the incident noise field, or the array geometry.

Some of the material in this section makes use of the work of Gaarder. (2)-(3)

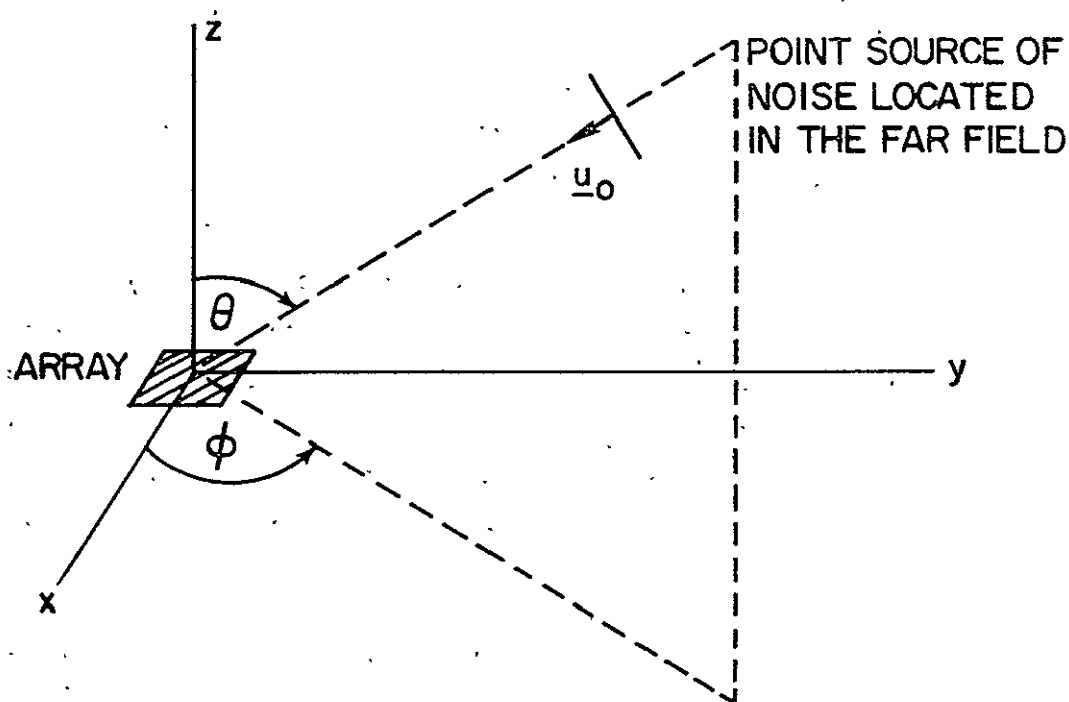


Fig. 2.3.1 Incident Noise Field

For simplicity, let us initially assume that the total incident noise field consists of one plane wave emanating from one source located on the surface of a sphere of infinite radius as shown in Fig 2.3.1. We will denote this plane wave by  $\underline{p}(\theta_o, \phi_o, \underline{x}, t)$  where  $\theta_o$  and  $\phi_o$  are spherical coordinates specifying the direction of propagation, which is also denoted by  $\underline{u}_o$ .

In complex notation

$$\underline{p}(\theta_o, \phi_o, \underline{x}, t) = \underline{p} e^{-j \underline{k} \cdot \underline{x}} e^{j \omega t} \quad (2.3.1)$$

where  $\underline{p} = \underline{p}(\theta_o, \phi_o)$  is a complex scalar random variable (for electromagnetic fields  $\underline{p}$  would have to be a complex vector random variable, but we are considering acoustic fields) and

$$\underline{k} = \text{wavenumber} = \frac{\omega}{c} \underline{u}_o(\theta_o, \phi_o)$$

$$\underline{x} = x \underline{x}_o + y \underline{y}_o + z \underline{z}_o$$

$$\underline{u}_o = -\sin \theta_o \cos \phi_o \underline{x}_o - \sin \theta_o \sin \phi_o \underline{y}_o - \cos \theta_o \underline{z}_o$$

An alternate way of writing  $p(\theta_o, \phi_o, \underline{x}, t)$  is

$$\begin{aligned} p(\theta_o, \phi_o, \underline{x}, t) &= p e^{+j(\omega t - \underline{k} \cdot \underline{x})} \\ &= p e^{j\omega(t - \underline{K} \cdot \underline{x})} \end{aligned} \quad (2.3.2)$$

where  $\underline{K} \equiv \frac{\underline{k}(\theta_o, \phi_o)}{\omega} = \frac{1}{c} \underline{u}_o(\theta_o, \phi_o)$  is independent of frequency. (2.3.3)

Since the actual noise sources we wish to investigate do not emit monochromatic waveforms but rather superpositions of monochromatic waveforms, let us change the assumption of one plane wave emanating from one source to an arbitrary superposition of plane waves emanating from one source.

In this case

$$\underline{k} = \underline{k}(\theta_o, \phi_o, \omega) = \frac{\omega}{c} \underline{u}_o(\theta_o, \phi_o)$$

$$\underline{K} = \frac{1}{c} \underline{u}_o(\theta_o, \phi_o) \text{ is still independent of frequency}$$

$$p(\theta_o, \phi_o, \underline{x}, t) = \int_{\omega} p(\theta_o, \phi_o, \omega) e^{j\omega(t - \underline{K} \cdot \underline{x})} d\omega \quad (2.3.4)$$

Noting that  $(t - \underline{K} \cdot \underline{x})$  is independent of frequency, we may define

$$p(\theta_o, \phi_o, \underline{x}, t) \equiv q(\theta_o, \phi_o, t - \underline{K} \cdot \underline{x}) \quad (2.3.5)$$

where  $q(\theta_o, \phi_o, t - \underline{K} \cdot \underline{x})$  for fixed  $\theta_o$  and  $\phi_o$  is a sample function of a stationary, zero-mean random process, with space-time covariance function

$$G_q(\theta_o, \phi_o, t_1 - t_2, \underline{x}_1 - \underline{x}_2) \equiv E \left\{ q^*(\theta_o, \phi_o, t_1 - \underline{K} \cdot \underline{x}_1) q(\theta_o, \phi_o, t_2 - \underline{K} \cdot \underline{x}_2) \right\} \quad (2.3.6)$$

Let us now drop the assumption of there being only one source located at coordinates  $(\theta_0, \phi_0)$  and instead assume that the noise field is generated by one point source on the infinite sphere corresponding to every different value of  $(\theta, \phi)$ . Thus the total noise field is given by

$$\underline{n}(\underline{x}, t) = \int_{\theta} \int_{\phi} \underline{q}(\theta, \phi, t - \underline{K} \cdot \underline{x}) d\Omega \quad (2.3.7)$$

We will assume that the sources are statistically independent of one another, implying that  $\underline{q}(\theta_1, \phi_1, t - \underline{K}(\theta_1, \phi_1) \cdot \underline{x})$  is independent of

$\underline{q}(\theta_2, \phi_2, t - \underline{K}(\theta_2, \phi_2) \cdot \underline{x})$  if  $(\theta_1, \phi_1) \neq (\theta_2, \phi_2)$ , i.e.

$$E \left\{ \underline{q}(\theta_1, \phi_1, t - \underline{K} \cdot \underline{x}_\alpha) \underline{q}(\theta_2, \phi_2, t - \underline{K} \cdot \underline{x}_\beta) \right\} = 0 \quad (2.3.8)$$

for  $(\theta_1, \phi_1) \neq (\theta_2, \phi_2)$

We may combine equations (2.3.6) and (2.3.8) to give

$$C_q(\theta_1, \phi_1, \theta_2, \phi_2, t_1 - t_2, \underline{x}_1 - \underline{x}_2) = C_q(\theta_1, \theta_2, t_1 - t_2, \underline{x}_1 - \underline{x}_2) \delta(\theta_1 - \theta_2, \phi_1 - \phi_2) \quad (2.3.9)$$

where

$$\int_0^\pi \int_0^{2\pi} \delta(\theta_1 - \theta_2, \phi_1 - \phi_2) \sin \theta_1 d\theta_1 d\phi_1 = 1 \quad (2.3.10)$$

Thus, the total noise field is stationary, with zero mean and space-time covariance

$$\begin{aligned} C_n(t_1 - t_2, \underline{x}_1 - \underline{x}_2) &= E \left\{ \underline{n}^*(\underline{x}_1, t_1) \underline{n}(\underline{x}_2, t_2) \right\} \\ &= \int_{\theta} \int_{\phi} C_q(\theta, \phi, t_1 - t_2, \underline{x}_1 - \underline{x}_2) d\Omega \end{aligned} \quad (2.3.11)$$

Note that, if the number of statistically independent noise sources is large, the resulting total noise field is gaussian, and the mean and covariance  $C_n$  completely describe the noise field.

Two simple special cases of the above general noise field (evaluated for the special case  $\underline{x}_1 = \underline{x}_2$  — we will later show that this is the only case we must consider explicitly, all other cases follow from this one by equation 2.3.16) are:

Monochromatic Noise

$$C_q(\theta, \phi, \tau, \underline{x}_1 - \underline{x}_2 = \underline{0}) = C_q(\theta, \phi, 0) e^{+j 2 \pi f_0 \tau} \quad (2.3.12a)$$

White Noise

$$C_q(\theta, \phi, \tau, \underline{x}_1 - \underline{x}_2 = \underline{0}) = C_q(\theta, \phi, 0) \delta(\tau) \quad (2.3.12b)$$

Let us now find the correlation between any two detector locations  $\underline{x}_1$  and  $\underline{x}_2$  in the x - y plane.

The noise incident upon a receiver located at  $\underline{x}_1$  is

$$n(\underline{x}_1, t) = \int_{\theta} \int_{\phi} q(\theta, \phi, t - \underline{K} \cdot \underline{x}_1) d\Omega \quad (2.3.13)$$

We will now let  $\underline{x}_1$  be the origin of our coordinate system, since only the magnitude and direction of the difference  $\underline{x}_2 - \underline{x}_1$  is of importance (this is because the noise sources are in the far field).

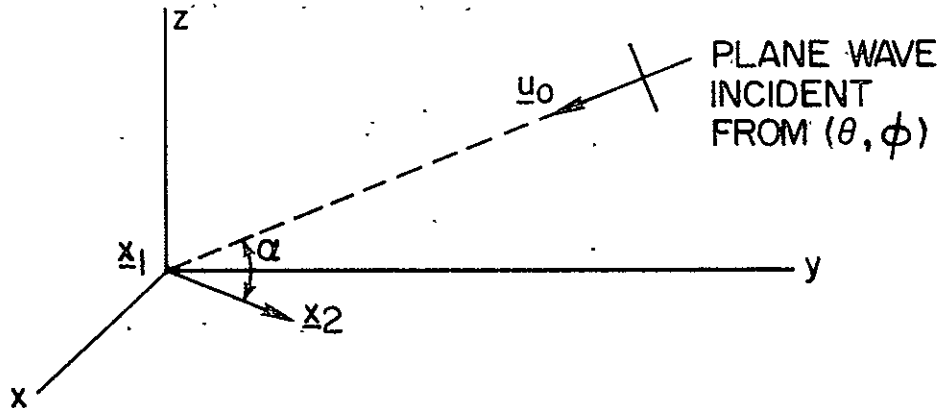


Fig. 2.3.2 Correlation between two detectors

We assume there is no attenuation as each plane wave comprising the noise field travels between the detectors at positions  $\underline{x}_1$  and  $\underline{x}_2$ . All plane waves, no matter what their frequency, move at the same velocity, because the medium is assumed to be homogeneous and isotropic.

Let  $\alpha$  be the angle between  $(\theta, \phi)$  and  $(\underline{x}_2 - \underline{x}_1)$ , i.e.  $\alpha = \alpha(\theta, \phi)$  is measured in the plane formed by the line  $\underline{x}_2 - \underline{x}_1$  and the direction of the incident plane wave  $\underline{u}_0$ . As we have the coordinates set up, with the noise field incident from the first octant and  $\underline{x}_2$  in the first quadrant, the noise wave hits  $\underline{x}_2$  before  $\underline{x}_1$  in time. Thus if the noise hits  $\underline{x}_1$  at time  $t$ , it hits  $\underline{x}_2$  at time  $t - T_{12} \cos \alpha$  where

$$T_{12} \cos \alpha \equiv - \frac{\underline{u}_0 \cdot (\underline{x}_2 - \underline{x}_1)}{C} \quad (2.3.14)$$

On the other hand, if the noise is at  $\underline{x}_2$  at time  $t$ , it is at  $\underline{x}_1$  at time  $t + T_{12} \cos \alpha$ .

Thus

$$\underline{n}(\underline{x}_1, t) = \underline{n}(\underline{x}_2, t - T_{12} \cos \alpha) \quad (2.3.15a)$$

$$\underline{n}(\underline{x}_2, t) = \underline{n}(\underline{x}_1, t + T_{12} \cos \alpha) \quad (2.3.15b)$$

The space-time correlation function of the noise process is

$$\begin{aligned} \phi_n(\tau, \underline{x}_1 - \underline{x}_2) &\equiv E \left\{ \underline{n}^*(\underline{x}_1, t) \underline{n}(\underline{x}_2, t - \tau) \right\} \\ &= E \left\{ \underline{n}^*(\underline{x}_1, t) \underline{n}(\underline{x}_1, t - \tau + T_{12} \cos \alpha) \right\} \\ &= C_n(\tau - T_{12} \cos \alpha, \underline{x}_1 - \underline{x}_1) \\ &= \int_{\theta} \int_{\phi} C_q(\theta, \phi, \tau - T_{12} \cos \alpha, \underline{0}) d\Omega \end{aligned} \quad (2.3.16)$$

Under the monochromatic noise assumption of equation (2.3.12a)

$$\phi_n(\tau, \underline{x}_1 - \underline{x}_2) = \int_{\theta} \int_{\phi} C_q(\theta, \phi, \underline{o}) e^{j 2 \pi f_o [\tau - T_{12} \cos \alpha]} d\Omega \quad (2.3.17a)$$

Under the white noise assumption of equation (2.3.12b)

$$\phi_n(\tau, \underline{x}_1 - \underline{x}_2) = \int_{\theta} \int_{\phi} C_q(\theta, \phi, \underline{o}) \delta[\tau - T_{12} \cos \alpha] d\Omega \quad (2.3.17b)$$

Equations (2.3.17) will be used in the multichannel filter point of view when we have to evaluate  $q_{ij} = E \{ n_i^*(t) n_j(t) \} = \phi_n(o, \underline{x}_i - \underline{x}_j)$ .

The total noise power incident at the origin (or at any detector) is given by  $\phi_n(o, \underline{o})$ . This follows by analogy with the power contained in a one dimensional random process whose autocorrelation function is  $R_x(\tau)$ , i. e. total power =  $\int_{\omega} S_x(\omega) d\omega = R_x(o)$ .

Noting that  $\underline{x}_1 - \underline{x}_2 = \underline{o}$  implies  $T_{12} \cos \alpha = o$ , we have, under both the monochromatic noise assumption and the white noise assumption

$$\phi_n(o, \underline{o}) = \int_{\theta} \int_{\phi} C_q(\theta, \phi, \underline{o}) d\Omega \quad (2.3.18)$$

Thus the spatial distribution of the noise power under either the monochromatic or white noise assumptions is

$$T(\theta, \phi) = C_q(\theta, \phi, \underline{o}) \quad (2.3.19)$$

In general, the equations we must use to transform between the detector pattern and multichannel filter viewpoints are, from equations (2.3.17) and (2.3.19):

Under the monochromatic noise assumption

$$\phi_n(\tau, \underline{x}_k - \underline{x}_l) = \int_{\theta} \int_{\phi} T(\theta, \phi) e^{j 2 \pi f_o [\tau - T_{kl} \cos \alpha]} d\Omega \quad (2.3.20a)$$

Under the white noise assumption

$$\phi_n(\tau, \underline{x}_k - \underline{x}_l) = \iint_{\theta, \phi} T(\theta, \phi) \delta[\tau - T_{kl} \cos \alpha] d\Omega \quad (2.3.20b)$$

Equations (2.3.20) are the results we have been striving for in this section. They express the space-time correlation functions  $\phi_n$  used in section 2.2 as direct functions of the incident noise power  $T(\theta, \phi)$  used in section 2.1.

We will now use these equations to show that under a monochromatic noise assumption (i. e. we will use equation (2.3.20a)), the detector pattern approach and the multichannel filter approach yield exactly the same values for the optimum SNR. We cannot show this is true for all possible spatial noise power distributions and all possible array configuration, because there is no general way of evaluating the integral in equation (2.3.20a). Because of this we will apply the theory developed above to three particular spatial noise power distributions and one particular array geometry. We will show, that under a monochromatic noise assumption, for these noise power distributions and this array configuration, the detector pattern approach and the multichannel filter approach yield exactly the same SNR results. Although we have used a particular array configuration and particular spatial noise power distributions, this was only done to simplify the evaluation of the integrals, and the equivalence can be seen to be independent of the incident spatial noise field and the array geometry.

The three spatial noise power distributions we will consider are:

1.  $T(\theta, \phi) = T(\theta, \phi) \delta(\theta - \theta_\beta, \phi - \phi_\beta)$
2.  $T(\theta, \phi) = T$  isotropic noise
3.  $T(\theta, \phi) = \begin{cases} T & \text{for } (\theta, \phi) \text{ in the first octant} \\ 0 & \text{otherwise} \end{cases}$

We will assume that the point detectors are equally spaced along the z axis, separated by a distance d.



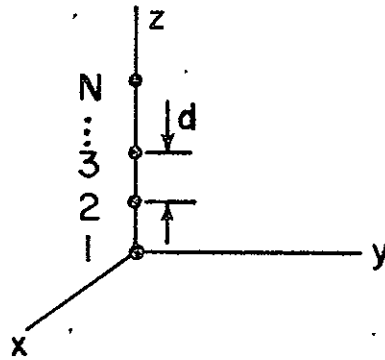


Fig. 2.3.3 Detector Array

In Appendix B we evaluate  $\phi_n(\tau, \underline{x}_k - \underline{x}_l)$  for the three spatial noise power distributions assuming the noise is temporally monochromatic and white. In Appendix C we evaluate the elements of the A matrix of section 2.1. In Appendix D we evaluate the elements of the Q matrix of section 2.2.

Using the results of the appendices, let us compare the results of sections 2.1 and 2.2. From section 2.1 we have as our expression for the optimum SNR achievable by using the detector pattern approach

$$\text{SNR} = \underline{V}_1^* A^{-1} \underline{V}_1$$

$$\text{where } \underline{V}_1^* = \begin{bmatrix} 1 & e^{-j 2 \pi \left(\frac{d}{\lambda}\right) \cos \theta_0} & \dots & e^{-j 2 \pi \left(\frac{d}{\lambda}\right) (N-1) \cos \theta_0} \end{bmatrix}$$

Note that we set  $\psi_i^0 = 2 \pi \left(\frac{d}{\lambda}\right) (i-1) \cos \theta_0$  because of our assumed array geometry.

Summarizing section 2.2, we have as our expression for the optimum SNR achievable by using the multichannel filter approach

$$\text{SNR} = \underline{U}_1^* Q^{-1} \underline{U}_1$$

$$\text{where } \underline{U}_1^* = \begin{bmatrix} 1 & e^{+j \omega \left[ -\frac{d}{c} \cos \theta_0 \right]} & \dots & e^{+j \omega \left[ -\frac{d}{c} (N-1) \cos \theta_0 \right]} \end{bmatrix}$$

Note that we set  $\tau_i = \frac{\underline{u}_0 \cdot \underline{r}_i}{c} = \frac{-d(i-1) \cos \theta_0}{c}$  because of our assumed array geometry.

$$\text{Since } \frac{\omega}{c} = \frac{2\pi f}{c} = \frac{2\pi f}{f\lambda} = \frac{2\pi}{\lambda}, \quad \underline{V}_1 \text{ and } \underline{U}_1 \text{ are equal.}$$

By comparing appendices C and D, we see that, for all three spatial noise fields considered, the A matrix of section 2.1 and the Q matrix of section 2.2 are equal, thus demonstrating that for monochromatic noise, we can optimize the SNR by using either the detector pattern or multichannel filter approach. Also note that from equations (2.1.14) and (2.2.16) the optimal current excitations and the optimal filter weights are equal, implying that the current excitations in the detector pattern approach correspond to the filter weights in the multichannel filter approach.

In conclusion, we have shown in this chapter, that under the monochromatic noise assumption, the detector pattern approach and the multichannel filter approach, are equivalent. Moreover, we saw that the current excitations of the detector pattern approach correspond to the filter weights of the multichannel filter approach. Again let us point out that although we have used a particular array configuration and particular spatial noise power distributions, this was only done to simplify the evaluation of certain integrals, and the equivalence can be seen to be independent of the array geometry and the incident spatial noise field.

In the next chapter, we will investigate the sensitivity of the SNR to small random changes in the detector locations and tap weights. We will have to use the equivalence developed in this chapter to derive an expression for this sensitivity. We will then show that when designing linear arrays where the spacing between detectors is less than one-half a wavelength, one should use tap weight values which maximize the SNR subject to a constraint on the above mentioned sensitivity, in order to keep this sensitivity within reasonable bounds.

# Appendix A

## Maximization of the SNR

$$\text{Maximize } L = \frac{\underline{I}^* C \underline{I}}{\underline{I}^* A \underline{I}} \quad \text{with respect to } \underline{I}.$$

Using the calculus of variations we get

$$\delta L = \frac{(\underline{I}^* A \underline{I}) \left[ (\delta \underline{I}^* C \underline{I}) + (\underline{I}^* C \delta \underline{I}) \right] - (\underline{I}^* C \underline{I}) \left[ (\delta \underline{I}^* A \underline{I}) + (\underline{I}^* A \delta \underline{I}) \right]}{(\underline{I}^* A \underline{I})^2} = 0$$

implying

$$\delta \underline{I}^* \left\{ C \underline{I} (\underline{I}^* A \underline{I}) - A \underline{I} (\underline{I}^* C \underline{I}) \right\} + \left\{ (\underline{I}^* A \underline{I}) \underline{I}^* C - (\underline{I}^* C \underline{I}) \underline{I}^* A \right\} \delta \underline{I} = 0$$

Since A and C are Hermitian

$$(\underline{I}^* A \delta \underline{I}) = (\delta \underline{I}^* A \underline{I})^*$$

$$(\underline{I}^* C \delta \underline{I}) = (\delta \underline{I}^* C \underline{I})^*$$

we have

$$\left\{ (\underline{I}^* A \underline{I}) \underline{I}^* C - (\underline{I}^* C \underline{I}) \underline{I}^* A \right\} \delta \underline{I} + \left[ \delta \underline{I}^* \left\{ C \underline{I} (\underline{I}^* A \underline{I}) - A \underline{I} (\underline{I}^* C \underline{I}) \right\} \right]^*$$

$$\text{Let } \underline{G} = C \underline{I} (\underline{I}^* A \underline{I}) - A \underline{I} (\underline{I}^* C \underline{I})$$

$$\text{thus } \delta \underline{I}^* \underline{G} + \left[ \delta \underline{I}^* \underline{G} \right]^* = 0$$

Since both of these terms are complex scalars and the second is the complex conjugate of the first, the real part of the complex scalar

must be zero, i. e.

$$\operatorname{Re} \left\{ \delta \underline{I}^* \underline{G} \right\} = 0$$

The only way this can be true for arbitrary  $\delta \underline{I}^*$  is if  $\underline{G} \equiv \underline{0}$ .

Thus

$$\underline{C} \underline{I} (\underline{I}^* \underline{A} \underline{I}) - \underline{A} \underline{I} (\underline{I}^* \underline{C} \underline{I}) = \underline{0}$$

By definition

$$\underline{C} = \underline{V}_1 \underline{V}_1^*$$

$$\underline{V}_1 (\underline{V}_1^* \underline{I}) (\underline{I}^* \underline{A} \underline{I}) - \underline{A} \underline{I} (\underline{I}^* \underline{V}_1) (\underline{V}_1^* \underline{I}) = \underline{0}$$

$$\underline{A} \underline{I} = \underline{V}_1 \frac{(\underline{I}^* \underline{A} \underline{I})}{(\underline{I}^* \underline{V}_1)}$$

$$\underline{I} = q \underline{A}^{-1} \underline{V}_1$$

Where the complex scalar  $q$  is given by

$$q \equiv \frac{(\underline{I}^* \underline{A} \underline{I})}{(\underline{I}^* \underline{V}_1)}$$

But the SNR is independent of the magnitude of  $\underline{I}$ , so when finding the value of  $\underline{I}$  which maximizes the SNR, we can drop the scalar  $q$ .

The direction of the optimum vector  $\underline{I}$ , which maximizes the SNR, is given by

$$\underline{I}_{\text{optimum}} = \underline{A}^{-1} \underline{V}_1$$

Appendix B Evaluation of  $\phi_n(\tau, \underline{x}_k - \underline{x}_l)$  for Temporally Monochromatic and White Noise.

Note that, for the array geometry of Fig 2.4.1, equation (2.3.14) becomes  $T_{kl} \cos \alpha \approx \frac{\cos \theta (\ell - k) d}{c}$  because  $\underline{u}_0(\theta, \phi) = -\sin \theta \cos \phi \underline{x}_0 - \sin \theta \sin \phi \underline{y}_0 - \cos \theta \underline{z}_0$  and  $(\underline{x}_l - \underline{x}_k) = (\ell - k) d \underline{z}_0$ .

If the noise is temporally monochromatic, for the three spatial noise power distributions under consideration, we have from equation (2.3.20a)

$$\text{case 1. } \phi_n(\tau, \underline{x}_k - \underline{x}_l) = T(\theta_\beta, \phi_\beta) e^{j 2 \pi f_o \left[ \tau - \frac{d(\ell - k)}{c} \cos \theta_\beta \right]}$$

$$\text{case 2. } \phi_n(\tau, \underline{x}_k - \underline{x}_l) = \int_0^\pi \int_0^{2\pi} T e^{j 2 \pi f_o \left[ \tau - \frac{d(\ell - k)}{c} \cos \theta \right]} \sin \theta d\theta d\phi$$

letting  $y = 2 \pi f_o \frac{d(\ell - k)}{c} \cos \theta$  and replacing  $\frac{c}{f_o}$  by  $\lambda$  gives

$$\phi_n(\tau, \underline{x}_k - \underline{x}_l) = \frac{2 T e^{j 2 \pi f_o \tau}}{\left( \frac{d}{\lambda} \right) (\ell - k)} \sin \left[ 2 \pi \left( \frac{d}{\lambda} \right) (\ell - k) \right]$$

$$\text{case 3. } \phi_n(\tau, \underline{x}_k - \underline{x}_l) = \int_0^{\pi/2} \int_0^{2\pi} T e^{j 2 \pi f_o \left[ \tau - \frac{d(\ell - k)}{c} \cos \theta \right]} \sin \theta d\theta d\phi$$

proceeding as in case 2, we get

$$\phi_n(\tau, \underline{x}_k - \underline{x}_l) = \frac{T \pi}{2} e^{j 2 \pi f_o \tau} e^{-j \pi \left( \frac{d}{\lambda} \right) (\ell - k)} \frac{\sin \left[ \pi \left( \frac{d}{\lambda} \right) (\ell - k) \right]}{\pi \left( \frac{d}{\lambda} \right) (\ell - k)}$$

If the noise is temporally white, for the three spatial noise power distributions under consideration, we have from equation (2.3.20b)

$$\text{case 1. } \phi_n(\tau, \underline{x}_k - \underline{x}_l) = T(\theta_\beta, \phi_\beta) \delta \left[ \tau - \frac{\cos \theta_\beta (\ell - k) d}{c} \right]$$

$$\text{case 2. } \phi_n(\tau, \underline{x}_k - \underline{x}_l) = T \int_0^\pi \int_0^{2\pi} \delta \left[ \tau - \frac{(\ell - k) d}{c} \cos \theta \right] \sin \theta d\theta d\phi$$

letting  $y = \frac{(\ell - k) d}{c} \cos \theta$  gives

$$\phi_n(\tau, \underline{x}_k - \underline{x}_l) = \begin{cases} \frac{-2\pi T c}{(\ell - k) d} & \text{if } |\tau| < \frac{(\ell - k) d}{c} \\ 0 & \text{otherwise} \end{cases}$$

$$\text{case 3. } \phi_n(\tau, \underline{x}_k - \underline{x}_l) = T \int_0^{\pi/2} \int_0^{\pi/2} \delta \left[ \tau - \frac{(\ell - k) d}{c} \cos \theta \right] \sin \theta d\theta d\phi$$

proceeding as in case 2, we get

$$\phi_n(\tau, \underline{x}_k - \underline{x}_l) = \begin{cases} \frac{-T c \pi}{2(\ell - k) d} & \text{if } 0 < \tau < \frac{(\ell - k) d}{c} \\ 0 & \text{otherwise} \end{cases}$$

## Appendix C Evaluation of the A matrix

From equation (2.1.12)

$$a_{kl} = \int_{\theta} \int_{\phi} e^{+j(\psi_k - \psi_l)} T(\theta, \phi) d\Omega$$

where

$$\psi_n = 2\pi \left[ \frac{x_n}{\lambda} \sin \theta \cos \phi + \frac{y_n}{\lambda} \sin \theta \sin \phi + \frac{z_n}{\lambda} \cos \theta \right] \text{ and}$$

$(x_n, y_n, z_n)$  is the position of the  $n^{\text{th}}$  detector.

For our array geometry the  $i^{\text{th}}$  detector is located on the  $z$  axis, at a distance  $z_i = d(i-1)$  from the origin, the above general expression becomes  $\psi_n = 2\pi \left(\frac{d}{\lambda}\right) (n-1) \cos \theta$ , thus

$$a_{kl} = \int_{\theta} \int_{\phi} e^{j 2\pi \left(\frac{d}{\lambda}\right) (k-l) \cos \theta} T(\theta, \phi) d\Omega$$

For the three spatial noise power distributions under consideration, we have

$$\text{case 1. } a_{kl} = e^{j 2\pi \left(\frac{d}{\lambda}\right) (k-l) \cos \theta_{\beta}} T(\theta_{\beta}, \phi_{\beta})$$

$$\begin{aligned} \text{case 2. } a_{kl} &= T \int_0^{\pi} \int_0^{2\pi} e^{j 2\pi \left(\frac{d}{\lambda}\right) (k-l) \cos \theta} \sin \theta d\theta d\phi \\ &= 4\pi T \frac{\sin \left[ 2\pi \left(\frac{d}{\lambda}\right) (k-l) \right]}{2\pi \left(\frac{d}{\lambda}\right) (k-l)} \end{aligned}$$

$$\text{case 3. } a_{kl} = T \int_0^{\pi/2} \int_0^{\pi/2} e^{j 2 \pi \left(\frac{d}{\lambda}\right) (k-l) \cos \theta} \sin \theta d\theta d\phi$$

$$= T \frac{\pi}{2} e^{j \pi \left(\frac{d}{\lambda}\right) (k-l)} \frac{\sin \left[ \pi \left(\frac{d}{\lambda}\right) (k-l) \right]}{\pi \left(\frac{d}{\lambda}\right) (k-l)}$$



## Appendix D      Evaluation of the Q matrix

From equations (2.2.4) and (2.3.16)

$$q_{k\ell} = E \left\{ n_k^* (t) n_\ell (t) \right\} = \phi_n (0, \underline{x}_k - \underline{x}_\ell)$$

In particular, when the noise is temporally monochromatic, for the three spatial noise power distributions under consideration, we have from Appendix B (remember  $\frac{f_0}{c} = \frac{1}{\lambda}$ )

$$\text{case 1.} \quad q_{k\ell} = \phi_n (0, \underline{x}_k - \underline{x}_\ell) = T(\theta_\beta, \phi_\beta) e^{-j 2\pi \left(\frac{d}{\lambda}\right) (\ell-k) \cos \theta_\beta}$$

$$\text{case 2.} \quad q_{k\ell} = \frac{4\pi T \sin \left[ 2\pi \left(\frac{d}{\lambda}\right) (\ell-k) \right]}{2\pi \left(\frac{d}{\lambda}\right) (\ell-k)}$$

$$\text{case 3.} \quad q_{k\ell} = \frac{T\pi e^{-j\pi \left(\frac{d}{\lambda}\right) (\ell-k)}}{2} \frac{\sin \left[ \pi \left(\frac{d}{\lambda}\right) (\ell-k) \right]}{\pi \left(\frac{d}{\lambda}\right) (\ell-k)}$$

## CHAPTER 3

### Error Analysis of Point Detector Arrays

If we were to design a point detector array or a multichannel filter to extract a signal, incident from direction  $(\theta_0, \phi_0)$ , from background noise, using the criterion of maximizing the SNR, as developed in chapter two, the following types of errors might affect the performance of our system:

1. Small random errors in the antenna excitations or filter coefficients (possibly due in part to round-off errors if we use a digital system to determine the filter coefficients).

2. Imperfect knowledge of the noise field.

That error type two is of importance is self-evident. However, the reader may ask if error type one is very important. It turns out that error type one can be of major importance as can be seen by considering the following problem:

Assume we wish to receive a signal propagating in the  $z$  direction, having wavenumber  $k_z = \frac{2\pi}{\lambda}$ , by using a linear array of  $N$  isotropic point detectors located along the  $z$  axis. Because of the sampling theorem, our first inclination would be to space the  $N$  detectors one-half wavelength apart ( $\frac{\lambda}{2}$ ), and then proceed to optimize the excitations so as to maximize the SNR. The question is, how much does error type one affect us if we use this spacing? It will be shown that for spacings between detectors of less than about one-half wavelength, the super-gain ratio, which is a measure of how much type one errors affect the detector pattern and thus the SNR, begins to get very large. This means that very small errors in the antenna excitations cause large variations in the received SNR. A better approach to use when the detectors are separated by less than a wavelength, would be to maximize the SNR subject to a constraint on the super-gain, or type one, error. This is one of the things we will investigate in this chapter.

Because the above mentioned types of errors are present in our system, the following questions arise:

1. If we optimize the processor as in chapter two, what are the effects of error type one on the SNR?

2. What is the optimum SNR we can achieve if we optimize the processor subject to a constraint on error type one?

3. Can we develop an adaptive algorithm which maximizes the SNR subject to a constraint on error type one?

The reason for undertaking this entire investigation is to answer question three - because the development of this type of algorithm will enable us to design array processors which will no longer significantly suffer from the deleterious effects of error types one and two that present day arrays suffer from.

In this chapter we will answer questions one and two. We will answer question three in chapters four, five and six.

### Section 3.1    Sensitivity of the SNR to Random Errors in the Detector Excitations and Locations.

Consider an array of  $N$  isotropic detectors placed at some prescribed positions in space whose Cartesian coordinates are given by  $\underline{x}_i$ ,  $i=1, \dots, N$ . Let  $(\theta_0, \phi_0)$  be the angular coordinates of the main beam, and  $I_i$  be the current excitation in the  $i^{\text{th}}$  detector. From equation (2.1.13) the SNR is given by

$$\text{SNR} = \frac{\underline{I}^* \underline{V}_1 \underline{V}_1^* \underline{I}}{\underline{I}^* \underline{\Lambda} \underline{I}} \quad (3.1.1)$$

where all quantities have been defined previously in section 2.1.

By the sensitivity of the SNR to random errors in the detector excitations and locations we mean the following: if we let the detector currents and positions be composed of a nominal term plus a random term, i.e.  $\underline{I} \rightarrow \underline{I}_n + \underline{I}_r$  and  $\underline{x}_i \rightarrow \underline{x}_{in} + \underline{x}_{ir}$ , the SNR is now defined as the expected value of equation (3.1.1). This expectation might turn out to be of the form

$$E \left\{ \frac{\underline{I}^* \underline{V}_1 \underline{V}_1^* \underline{I}}{\underline{I}^* \underline{\Lambda} \underline{I}} \right\} = \frac{\underline{I}_n^* \underline{V}_1 \underline{V}_1^* \underline{I}_n}{\underline{I}_n^* \underline{\Lambda} \underline{I}_n} + \text{an additional term, and we would}$$

then define the ratio of the additional term to the nominal term as our sensitivity factor. The calculation of this expected value, as it stands, is exceedingly complex. However, the SNR in equation (3.1.1) may also be expressed as

$$\text{SNR} = \frac{\Phi(\underline{u}_0)}{\int_{\Omega} \Phi(\underline{u}) T(\underline{u}) d\Omega} \quad (3.1.2)$$

where  $\Phi(\underline{u})$  is the array power pattern  $\Phi(\underline{u}) = |\underline{I}^* \underline{V}|^2$ ,  $\Phi(\underline{u}_0) = |\underline{I}^* \underline{V}_1|^2$  is the value of the power pattern at  $(\theta_0, \phi_0)$ , and  $T(\underline{u})$  is the incident noise power. Again, if we let the detector currents and locations be random, the calculation of the expected value of equation (3.1.2) is exceedingly complex.

However, equation (3.1.2) indicates to us that we can use the super-gain ratio, which is a measure of the sensitivity of the power pattern  $\Phi(\underline{u})$  to random errors in the detector excitations and positions, as an alternate measure of the sensitivity of the SNR to random errors in the detector excitations and positions.

An intuitive justification for this is as follows:  $I(\underline{u})$  is the power pattern. Since the signal is incident from direction  $\underline{u}_0$ , the power pattern is usually designed so as to peak up in the  $\underline{u}_0$  direction, e.g.

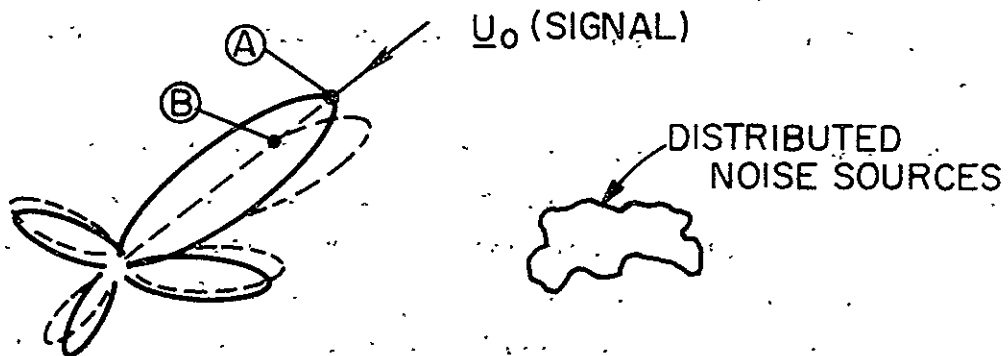


Fig. 3.1.1 Typical Power Pattern

The solid line in Fig 3.1.1 represents the theoretical power pattern while the dashed line represents the actual pattern we may get due to random

errors in current excitations and detector locations. Small changes in the power pattern affect the numerator of the SNR much more than the denominator because the numerator is proportional to the pattern while the denominator is proportional to the integral of the power pattern over all space, which doesn't change as much. Put another way, if the power pattern changes slightly, the main reason for the change in the SNR is because the signal power received by the array drops from level A to level B. While the noise power received by the array changes, it does not change to as great an extent as did the signal power received. Thus our premise is that

$$\Delta [ \text{SNR} ] \propto \Delta [ I(\underline{u}) ]$$

The super-gain ratio Q is derived in Appendix A and is given by equation (A15)

$$Q \equiv \frac{\underline{I}^* \underline{I}}{\int_{\Omega} \underline{I}^* \underline{V} \underline{V}^* \underline{I} \, d\Omega} = \frac{\underline{I}^* \underline{I}}{\underline{I}^* \underline{B} \underline{I}} \quad (3.1.3)$$

where  $\underline{B} \equiv \int_{\Omega} \underline{V} \underline{V}^* \, d\Omega$

Q is a function of the spacing between detectors through  $\underline{V}$ , and through  $\underline{I}$  is also a function of the signal location (or main beam direction) and the noise field (i. e. assuming we use that value of  $\underline{I}$  which maximizes the SNR).

To investigate how the SNR and Q factor behave as a function of array geometry, we shall focus on the special case of Fig 3.1.2, consisting of a linear array of four isotropic detectors embedded in a uniform noise field (i. e.  $T(\theta, \phi) = 1$  for  $0 \leq \theta \leq \pi$ ,  $0 \leq \phi < 2\pi$ ), whose main beam is at broadside ( $\theta_0 = 0$ ) or endfire ( $\theta_0 = \frac{\pi}{2}$ ,  $\phi_0 = 0$ ), and whose current excitation is given by the optimum value of  $\underline{I}$  we found in chapter two (i. e. that value of  $\underline{I}$  which maximizes the SNR).

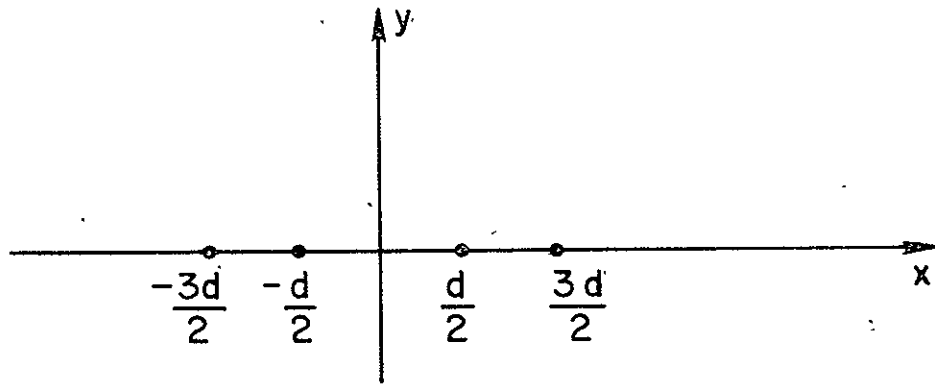


Fig.3.1.2 Four element linear array

Before we can obtain numerical results, we need the explicit form of the  $A$  matrix in the SNR expression for the case where  $T(\theta, \phi) = 1$  for all values of  $\theta$  and  $\phi$ , and of the matrix  $\int_{\Omega} \underline{V} \underline{V}^* d\Omega$  in the  $Q$  factor expression. Because of our choice of an isotropic noise field, these matrices become identical, and, in this case, the elements of  $A$ , denoted by  $a_{kl}$ , can be integrated out in closed form for planar arrays of isotropic elements. Assuming the detectors are in  $xy$  plane, the elements of  $A$  are given by

$$\begin{aligned}
 a_{kl} &= a_{lk}^* = \int_0^{2\pi} \int_0^\pi e^{+j\psi_k} e^{-j\psi_l} \sin\theta d\theta d\phi \\
 &= \int_0^{2\pi} \int_0^\pi e^{-j2\pi \left[ \left( \frac{x_k - x_l}{\lambda} \right) \sin\theta \cos\phi + \left( \frac{y_k - y_l}{\lambda} \right) \sin\theta \sin\phi \right]} \sin\theta d\theta d\phi
 \end{aligned} \tag{3.1.4}$$

We may rewrite the integrand by noting the following identity

$$A_1 \cos\phi + A_2 \sin\phi = \operatorname{Re} \left[ A_1 e^{j\phi} + A_2 e^{j(\phi - \pi/2)} \right]$$

$$A_1 e^{j\phi} + A_2 e^{j\phi} e^{-j\pi/2} = A_1 e^{j\phi} - j A_2 e^{j\phi} = [A_1 + j(-A_2)] e^{j\phi}$$

$$= \sqrt{A_1^2 + A_2^2} e^{j \tan^{-1} \frac{-A_2}{A_1}} e^{j\phi} = \sqrt{A_1^2 + A_2^2} e^{j(\phi - \tan^{-1} \frac{A_2}{A_1})}$$

Taking the real part gives the result

$$A_1 \cos \phi + A_2 \sin \phi = \sqrt{A_1^2 + A_2^2} \cos(\phi - \tan^{-1} \frac{A_2}{A_1})$$

thus

$$a_{kl} = \int_0^{2\pi} \int_0^\pi e^{j 2\pi \sin \theta \rho_{kl} \cos(\phi - \gamma_{kl})} \sin \theta d\theta d\phi \quad (3.1.5)$$

where

$$\rho_{kl} \equiv \sqrt{\left(\frac{x_k - x_l}{\lambda}\right)^2 + \left(\frac{y_k - y_l}{\lambda}\right)^2} \quad (3.1.6)$$

$$\gamma_{kl} \equiv \tan^{-1} \frac{y_k - y_l}{x_k - x_l} \quad 0 \leq \gamma_{kl} < \pi \quad (3.1.7)$$

Note that  $\lambda_{kl}$  is a multivalued function, and since it appears in the integrand, it must be restricted. We will restrict  $\lambda_{kl}$  to the range  $0 \leq \lambda_{kl} < \pi$ . However, when we do this, if  $\lambda_{kl}$  appears explicitly in the resulting formula we get for  $a_{kl}$  we can not use the formula to calculate both  $a_{kl}$  and  $a_{lk}$  because we will not satisfy the requirement that  $a_{kl} = a_{lk}^*$  due to the restriction on  $\gamma$ . The procedure to use is as follows: If  $\gamma$  appears in the formula for  $a_{kl}$ , use the formula to evaluate  $a_{kl}$  for  $k$  strictly less than  $l$ , and evaluate  $a_{kl}$  for  $k > l$  by computing  $a_{lk}^*$ . If  $\gamma$  does not appear in the formula for  $a_{kl}$  (this is the result we will obtain in our problem, but we get this only because of the particular way we defined  $\underline{V}$  and  $\underline{I}$ ), there is no problem. In either case, to evaluate  $a_{kk}$ ,  $\gamma_{kk}$  is indeterminate and hence we must evaluate the diagonal terms separately.

Since 
$$\frac{1}{2\pi} \int_0^{2\pi} e^{j x \cos (\phi - \gamma_{k\ell})} d\phi = J_0(x)$$

$$a_{k\ell} = 2\pi \int_0^\pi \sin \theta J_0(2\pi \rho_{k\ell} \sin \theta) d\theta \quad (3.1.8)$$

But 
$$\int_0^\pi J_0(x \sin \theta) \sin \theta d\theta = 2 \frac{\sin x}{x}$$

$$a_{k\ell} = 4\pi \left[ \frac{\sin(2\pi \rho_{k\ell})}{2\pi \rho_{k\ell}} \right] \quad \text{for } k \neq \ell \quad (3.1.9)$$

If  $k = \ell$  we have

$$a_{kk} = \int_0^{2\pi} \int_0^\pi \sin \theta d\theta d\phi = 4\pi \quad (3.1.10)$$

For the special case of the four element linear array shown in Fig 3.1.2, the elements of the A matrix are given by

$$A = \begin{bmatrix} 4\pi & 2\frac{\lambda}{d} \sin 2\pi \frac{d}{\lambda} & \frac{\lambda}{d} \sin 4\pi \frac{d}{\lambda} & \frac{2}{3} \frac{\lambda}{d} \sin 6\pi \frac{d}{\lambda} \\ 2\frac{\lambda}{d} \sin 2\pi \frac{d}{\lambda} & 4\pi & 2\frac{\lambda}{d} \sin 2\pi \frac{d}{\lambda} & \frac{\lambda}{d} \sin 4\pi \frac{d}{\lambda} \\ \frac{\lambda}{d} \sin 4\pi \frac{d}{\lambda} & 2\frac{\lambda}{d} \sin 2\pi \frac{d}{\lambda} & 4\pi & 2\frac{\lambda}{d} \sin 2\pi \frac{d}{\lambda} \\ \frac{2}{3} \frac{\lambda}{d} \sin 6\pi \frac{d}{\lambda} & \frac{\lambda}{d} \sin 4\pi \frac{d}{\lambda} & 2\frac{\lambda}{d} \sin 2\pi \frac{d}{\lambda} & 4\pi \end{bmatrix} \quad (3.1.11)$$



The optimum (with respect to maximum SNR) value of  $\underline{I}$  is given by equation (2.1.18)

$$\underline{I}_{\text{opt}} = A^{-1} \underline{V}_1$$

Using this value of  $\underline{I}$ , we found in chapter 2 that the SNR is given by

$$\text{SNR} = \underline{V}_1^* A^{-1} \underline{V}_1 \quad (3.1.12)$$

Again using this value of  $\underline{I}$ , the Q factor is given by

$$Q = \frac{\underline{I}^* \underline{I}}{\underline{I}^* A \underline{I}} = \frac{\underline{V}_1^* [A^{-1}]^2 \underline{V}_1}{\underline{V}_1^* A^{-1} \underline{V}_1} \quad (3.1.13)$$

If the main beam is at broadside ( $\theta_o = 0$ ) then, in our example

$$\underline{V}_1 = \begin{bmatrix} e^{j\psi_1^o} \\ e^{j\psi_2^o} \\ e^{j\psi_3^o} \\ e^{j\psi_4^o} \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad (3.1.14)$$

If the main beam is at endfire ( $\theta_o = \frac{\pi}{2}$ ,  $\phi_o = 0$ ) then, in our example

$$\underline{V}_1 = \begin{bmatrix} e^{j\psi_1^o} \\ e^{j\psi_2^o} \\ e^{j\psi_3^o} \\ e^{j\psi_4^o} \end{bmatrix} = \begin{bmatrix} e^{j(-3\pi \frac{d}{\lambda})} \\ e^{j(-\pi \frac{d}{\lambda})} \\ e^{j(\pi \frac{d}{\lambda})} \\ e^{j(3\pi \frac{d}{\lambda})} \end{bmatrix} \quad (3.1.15)$$

Similar results can be obtained for the ten element linear array shown below in Fig 3.1.3

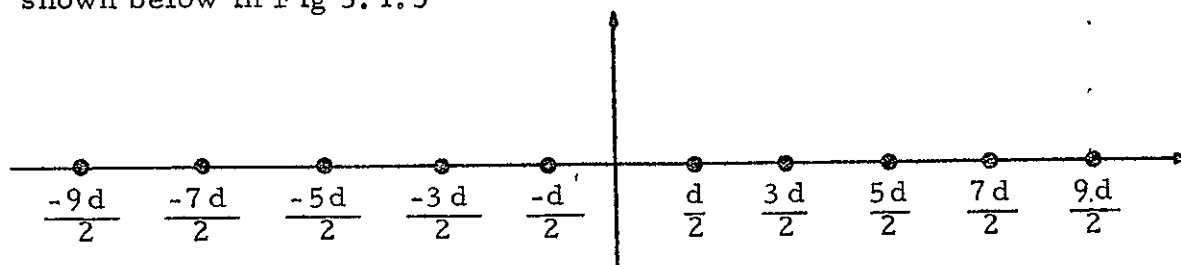


Fig 3.1.3 Ten element linear array

The following graphs of SNR and  $Q$  vs  $\frac{d}{\lambda}$  were obtained for four and ten element linear arrays, in isotropic noise, when the main beam was at broadside and endfire, using the optimum excitation:

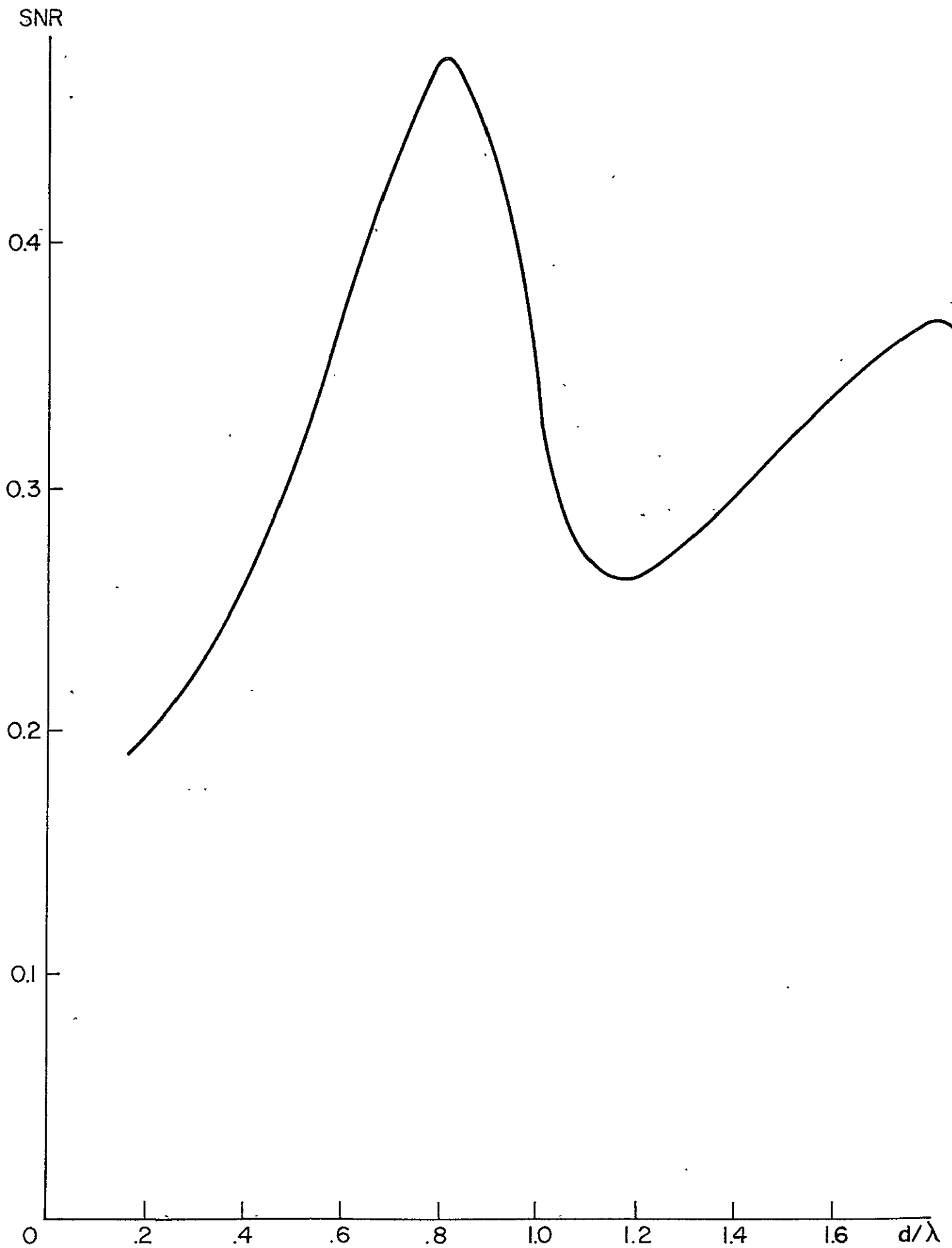


Fig. 3.1.4 Four Element Array:- Broadside Signal

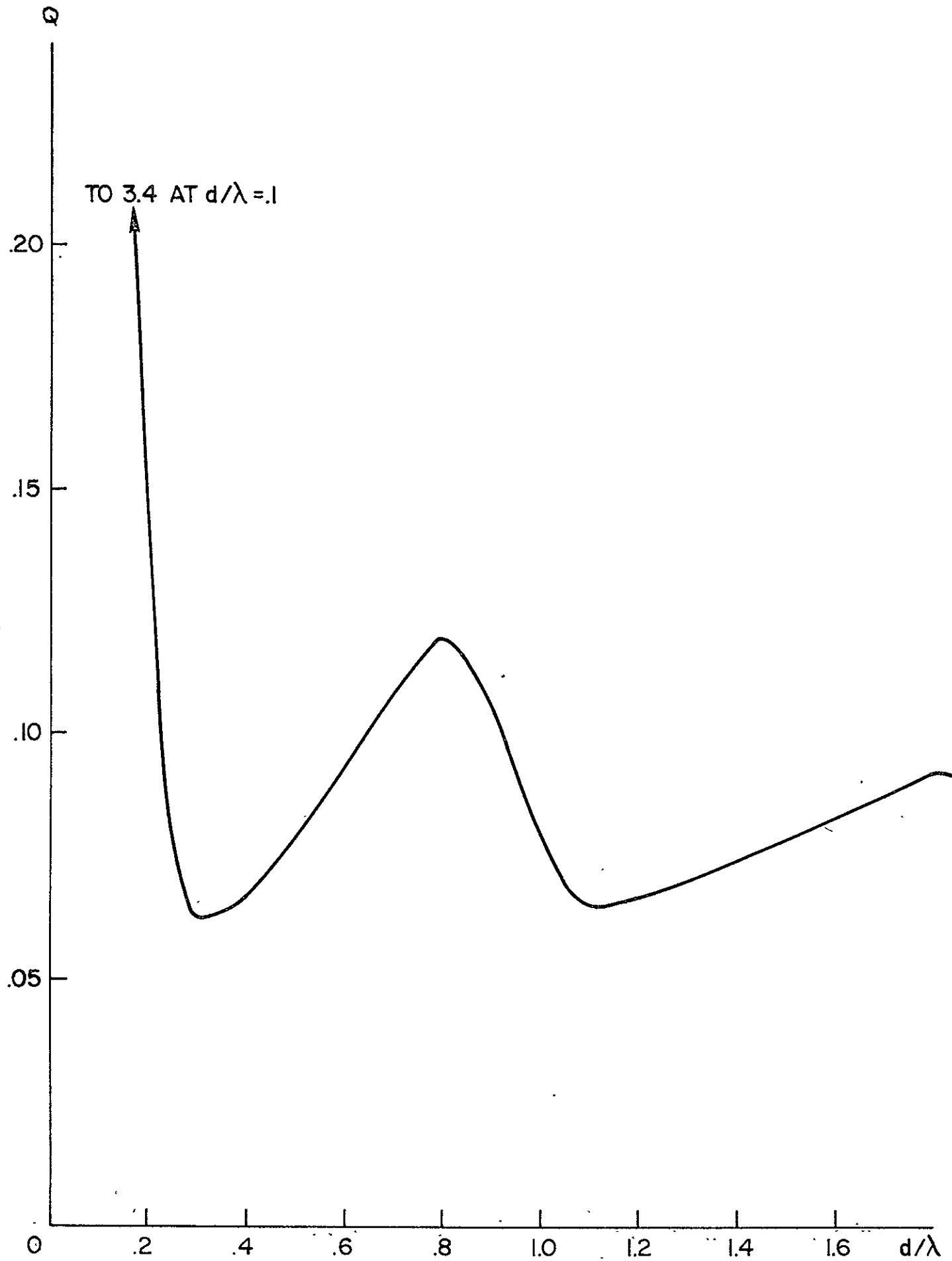


Fig. 3.1.5 Four Element Array - Broadside Signal

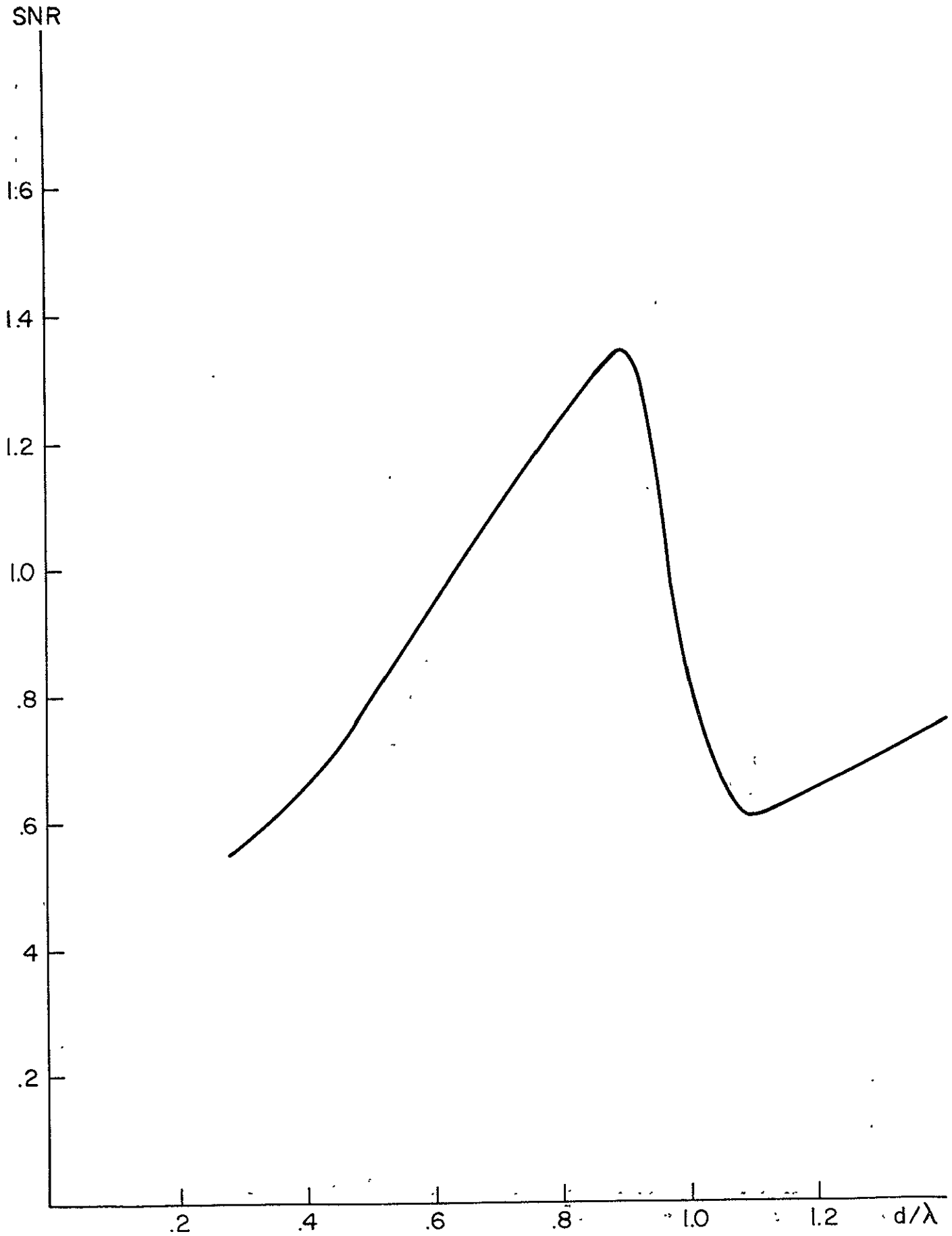


Fig. 3.1.6 Ten Element Array - Broadside Signal

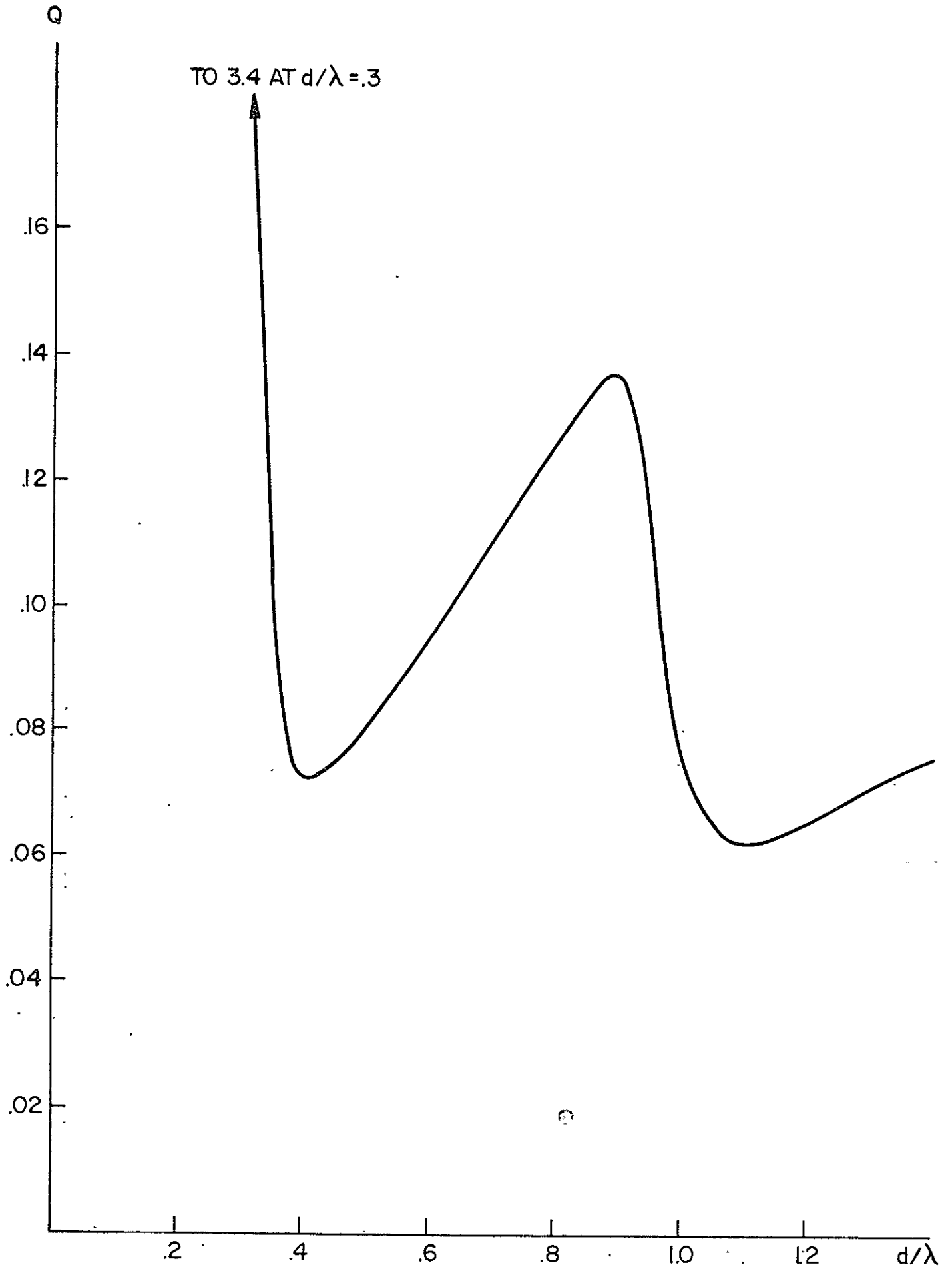


Fig. 3.1.7 Ten Element Array - Broadside Signal

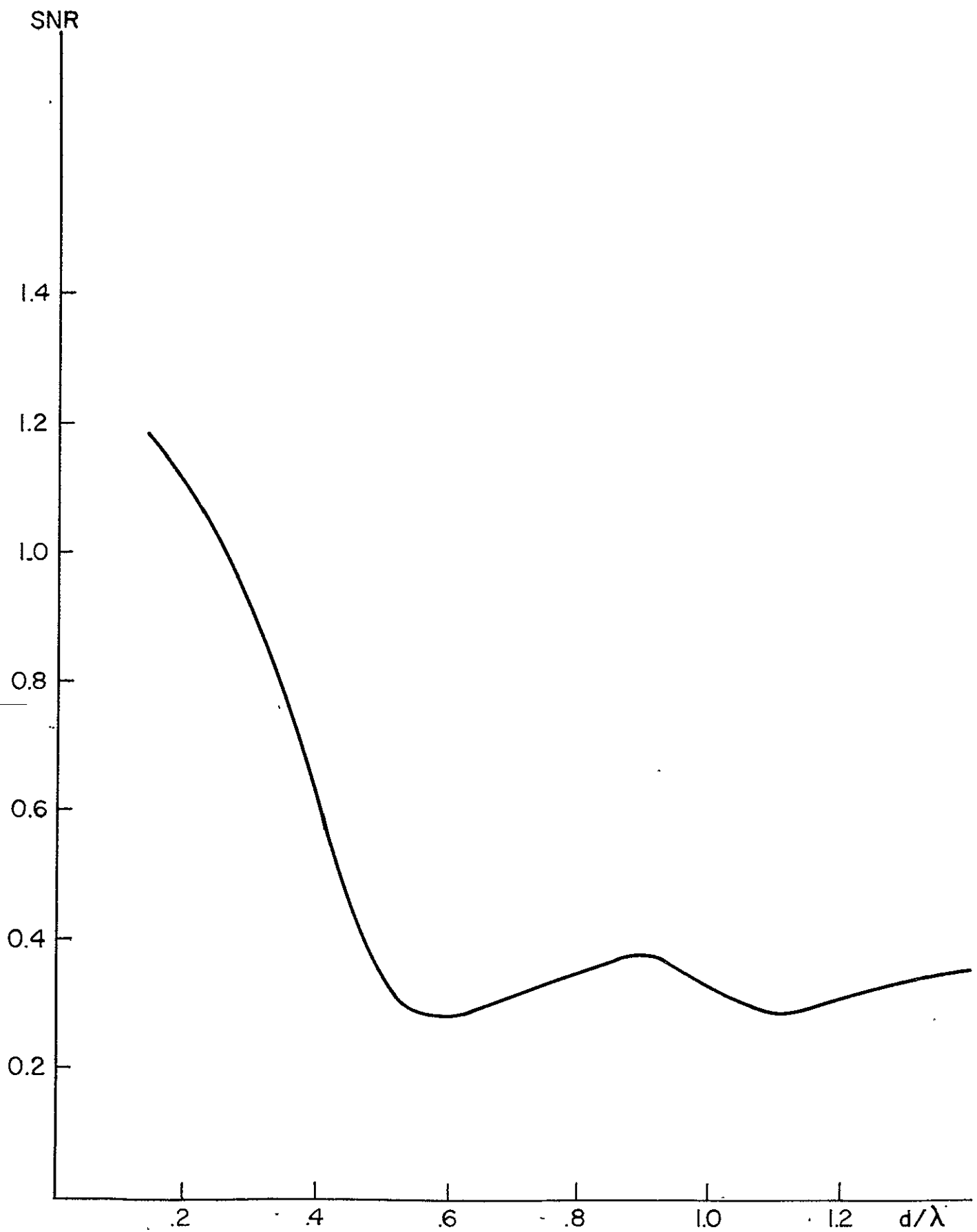


Fig. 3.1.8 Four Element Array - Endfire Signal

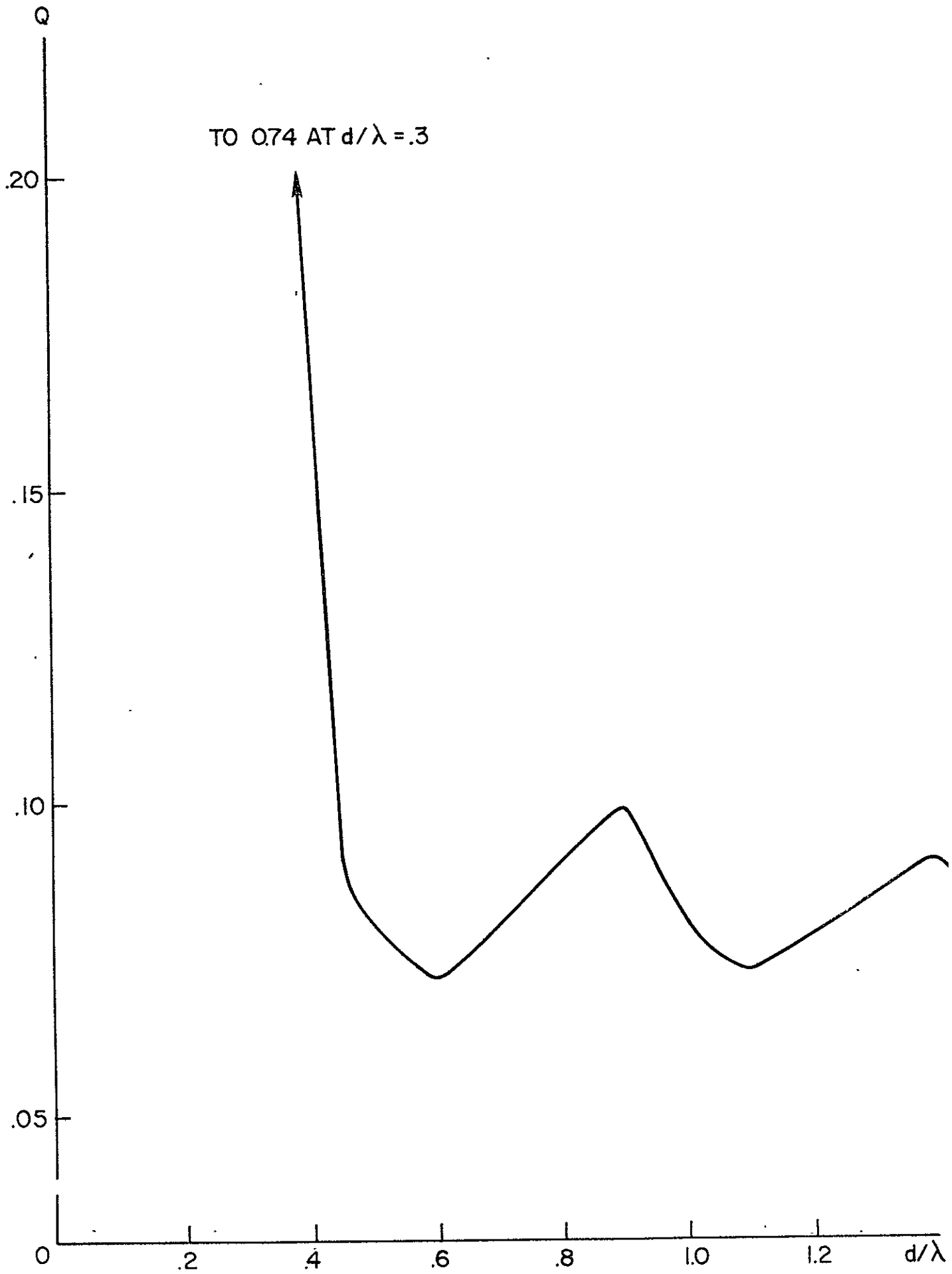


Fig. 3.1.9 Four Element Array - Endfire Signal



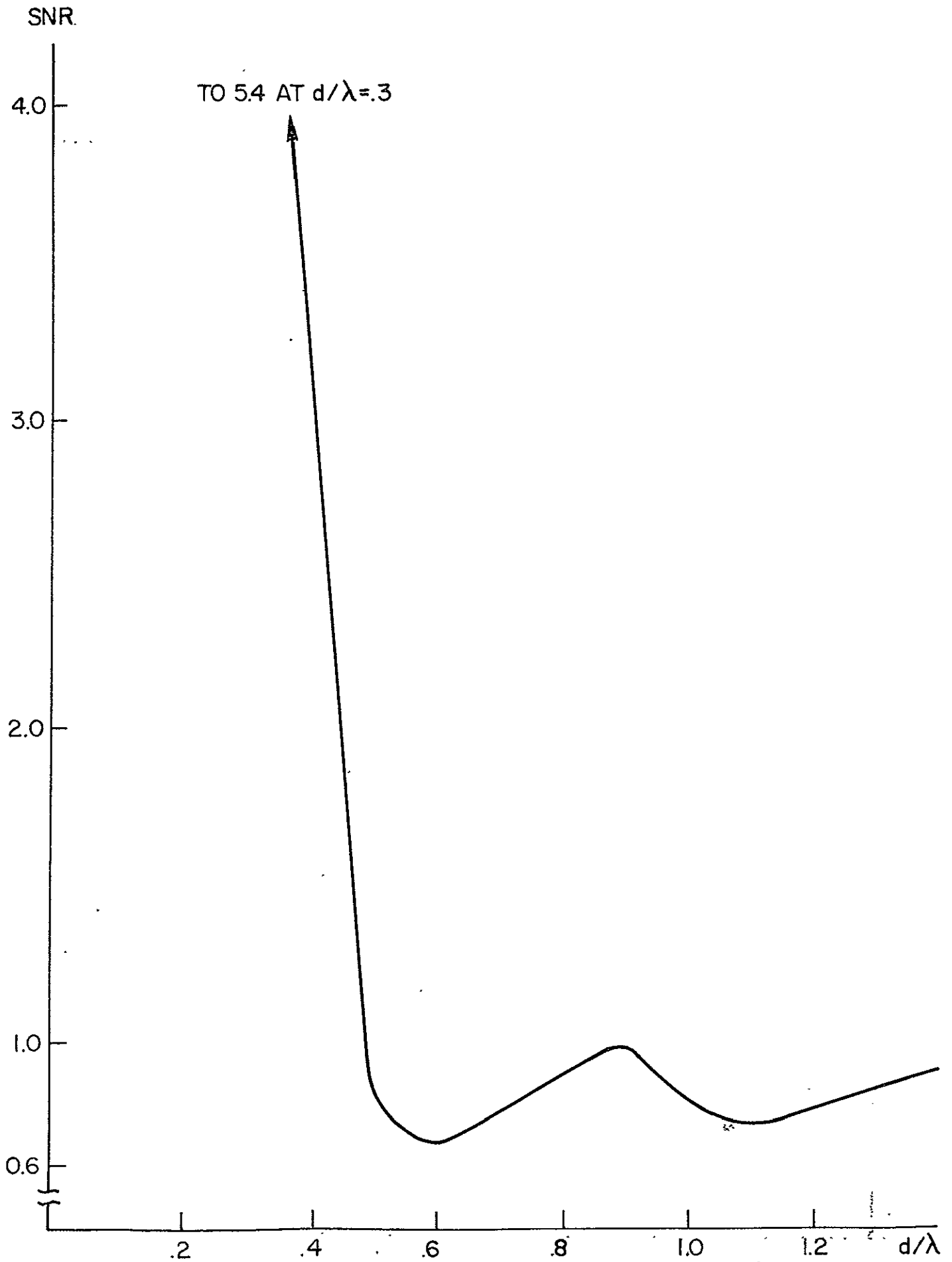


Fig. 3.1.10 Ten Element Array - Endfire Signal

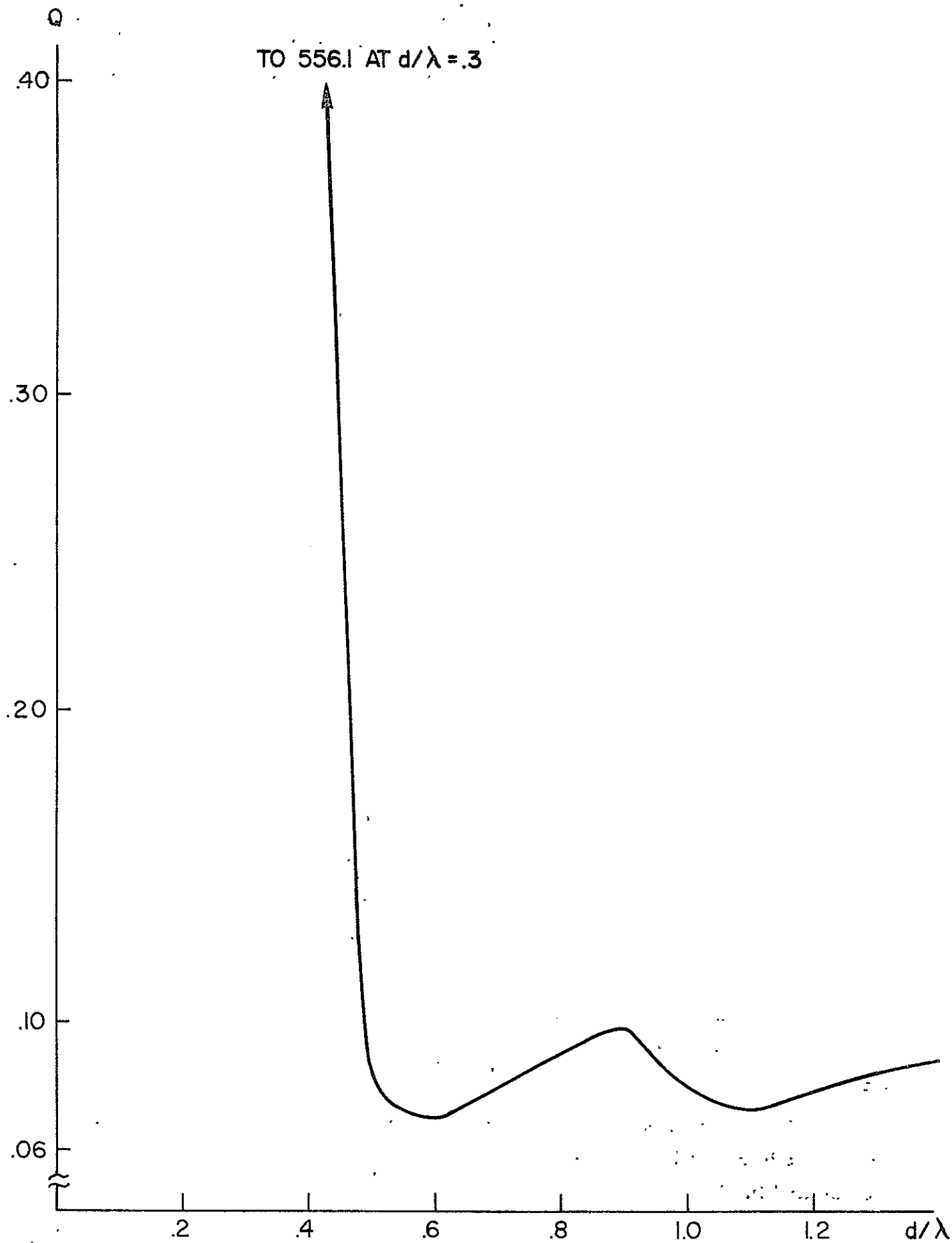


Fig. 3.1.11 Ten Element Array - Endfire Signal

By comparing Figs. 3.1.4 and 3.1.6, 3.1.5 and 3.1.7, 3.1.8 and 3.1.10, 3.1.9 and 3.1.11 we see that the general shape of the curves and the ratios of the maxima to the minima of each curve is independent of the number of elements (four vs ten) in the array. Hence in our future work we will only consider four element arrays in order to conserve computer time.

With reference to Figs. 3.1.4 and 3.1.5 notice that if we use those current excitations which maximize the SNR, the SNR and Q factor that we will get when the signal impinges from broadside can vary between 0.2 and 0.5 (a ratio of 1:2.5) and 0.05 to 0.15 (a ratio of 1:3) respectively, depending upon what spacing we use between detectors as long as it is greater than  $0.2\lambda$ .

Aside: Note that the graphs only cover the region up to  $d = 1.8\lambda$  because this is the region of interest to us; however, if we extended, for example, Fig 3.1.4, it looks as follows

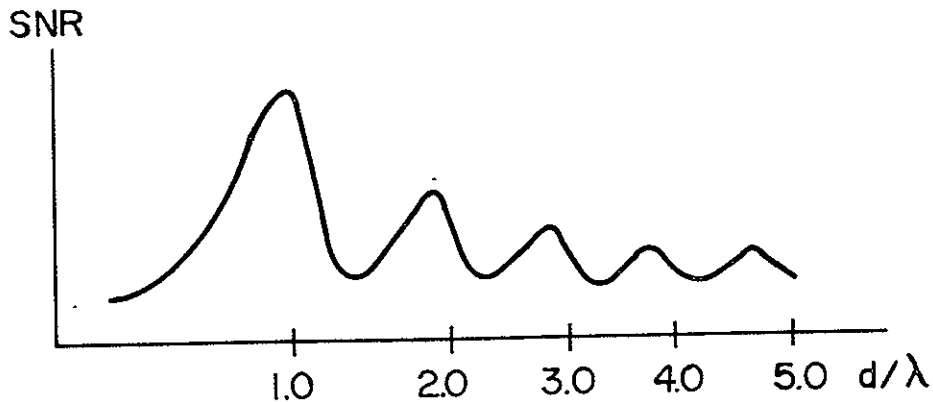


Fig. 3.1.12 Extension of Fig. 3.1.4

and all the other graphs behave similarly. Note also that our graphs don't cover the region  $d = 0$  to  $d = 0.2\lambda$  because in this region, mutual coupling effects between detectors come into play, and our analysis does not take this into account.

This means that for this array geometry, when the signal impinges from broadside, it is relatively unimportant what spacing between detectors we use and furthermore, it is acceptable for us to design the array (i.e.

choose the current excitations or tap weights) by maximizing the SNR alone - rather than designing the array by maximizing the SNR subject to a constraint on the Q factor - because the Q factor which results from the use of the first design procedure will never be excessive.

However, with reference to Figs. 3.1.8 and 3.1.9 notice that if we use those current excitations which maximize the SNR, the SNR and Q factor we will get when the signal impinges from endfire can vary between 0.2 and 1.0 (a ratio of 1:5) and 0.06 to a number well exceeding 0.74 (a ratio very much greater than 1:12) respectively, depending upon what spacing we use between detectors as long as it is greater than  $0.2 \lambda$ . This means that for this same array geometry, when the signal impinges from endfire, the spacing between detectors that we use is relatively important, i.e. we would prefer to space the detectors as close together as possible, however if we do this, the Q factor, which is a measure of the sensitivity of the SNR to the random fluctuations in the tap weights will be so large as to make the array processor useless.

The conclusion we draw from these graphs is that if we are going to use a certain detector array and we are not sure a priori that for all possible incident signal directions the Q factor never gets too large when we use those current excitations (or tap weights) which maximize the SNR, we must instead use those excitations which maximize the SNR (equation 3.1.1) subject to a constraint on the super-gain ratio (equation 3.1.3). We will see how to find these excitations in the next section.

### Section 3.2 Maximization of the SNR subject to a constraint on the super-gain ratio.

The problem is to maximize  $\frac{\underline{I}^* \underline{V}_1 \underline{V}_1^* \underline{I}}{\underline{I}^* \underline{A} \underline{I}}$  subject to the constraint

$$Q = \frac{\underline{I}^* \underline{I}}{\underline{I}^* \underline{B} \underline{I}} \quad (19)$$

Appendix B summarizes the work of Lo, Lee, and Lee recently developed a numerical technique of solving this problem. However, their work yields a (sometimes complex) polynomial equation whose roots (when found numerically) can then be used to calculate the value of  $\underline{I}$  which is the solution to the problem. Our contribution makes use of a state variable technique which enables us to reduce L, Lee and Lee's numerical problem from one of finding the complex roots of a high order polynomial with complex coefficients (in all the specific numerical cases treated in their paper the coefficients of the polynomials were real, but this is not necessarily true in general and is not true in the second example we will consider in this section) to one of finding the eigenvalues of a real matrix, which is considerably faster to do.

Since we can only get numerical results for particular examples, we will consider the following two specific problems:

1. Solve for that value of  $\underline{I}$  which will maximize the SNR subject to the constraint  $Q = .08$  for a linear array of four isotropic detectors spaced  $d = 0.8\lambda$  apart, embedded in a uniform noise field ( $T(\theta, \phi) = 1$  for  $0 \leq \theta \leq \pi$ ,  $0 \leq \phi \leq 2\pi$ ), whose main beam is at broadside ( $\theta_0 = 0$ ). From Fig 3.1.5 we see that if we did not constrain  $Q$ , but instead used that value of  $\underline{I}$  which maximized the SNR, we would get a value of  $Q$  equal to approximately 0.12.

2. Solve for that value of  $\underline{I}$  which will maximize the SNR subject to the constraint  $Q = .11$  for a linear array of four isotropic detectors spaced  $d = 0.4\lambda$  apart, embedded in a uniform noise field whose main beam is at endfire ( $\theta_0 = \pi/2$ ,  $\phi_0 = 0$ ). From Fig 3.1.9 we see that if we did not constrain  $Q$ , but instead used that value of  $\underline{I}$  which maximized the SNR, we would get a value of  $Q$  equal to approximately 0.18.

We will use Lo, Lee and Lee's method to do the first example, and our method to do the second. As far as the first example is concerned,  $\underline{V}_1 = \text{col } [1 \ 1 \ 1 \ 1]$  and we may

choose for our complete set (see Appendix B) the following vectors:

$$\underline{a}_1 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad \underline{a}_2 = \begin{bmatrix} -1 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad \underline{a}_3 = \begin{bmatrix} -1 \\ 0 \\ 1 \\ 0 \end{bmatrix} \quad \underline{a}_4 = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \end{bmatrix} \quad (3.2.1)$$

The W matrix (equation B8) has vectors  $\underline{a}_1, \underline{W}_2, \underline{W}_3, \underline{W}_4$  as columns, where

$$\underline{W}_i = (0.08 A - I) \underline{a}_i s^2 + 2 A \underline{a}_i s + A (0.08 A - I)^{-1} A \underline{a}_i; i = 2, 3, 4 \quad (3.2.2)$$

The elements of this matrix are real polynomials in  $s$  of degree two, except for the first column whose elements are all equal to one. Setting the determinant of this W matrix equal to zero results in a polynomial of sixth degree in  $s$  being equal to zero. After solving for the six roots, we take the real roots (since we know  $s$  is real) and substitute them into equation (B5) to determine the possible values of  $\underline{I}$ , i. e.

$$\underline{I} = [A - s I + 0.08 s A]^{-1} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad (3.2.3)$$

We now take these values of  $\underline{I}$  and substitute them into the expressions for Q and SNR. The solution we are looking for is given by the  $\underline{I}$  which satisfies

$$Q = \frac{\underline{I}^* \underline{I}}{\underline{I}^* A \underline{I}} = 0.08 \text{ and gives the highest value of the SNR } = \frac{\underline{I}^* \underline{V}_1 \underline{V}_1^* \underline{I}}{\underline{I}^* A \underline{I}}.$$

Numerically, we found the following six roots of the polynomial, the real roots being allowable values of  $s$ ; corresponding to these four allowable values of  $s$  we found the values of the Q factor, corresponding to the two values of  $s$  for which the Q factor is equal to 0.08 we found the two values of the SNR.

| S               | Q     | SNR   |
|-----------------|-------|-------|
| 121.0 + j 0.198 | ----  | ----  |
| 121.0 - j 0.198 | ----  | ----  |
| -112.7          | 0.080 | 0.058 |
| -52.2           | 0.080 | 0.187 |
| -61.8           | 0.070 | 0.084 |
| -61.1           | 0.071 | 0.090 |

The solution to the first problem, i. e. that value of  $\underline{I}$  which maximizes the SNR subject to the constraint  $Q = 0.08$  for a broadside array is the value of  $\underline{I}$  corresponding to  $s = -52.160$ . For this value of  $s$ ,  $\underline{I}$  is given by

$$\underline{I} = \begin{bmatrix} 0.086 \\ 0.007 \\ 0.007 \\ 0.086 \end{bmatrix}$$

and the maximum SNR we can achieve subject to the constraint  $Q = 0.08$  is  $SNR = 0.187$ .

The second example is more complicated, because the vector space we are working in consists of complex vectors (e.g.  $\underline{a}_1$ ) over a complex scalar field (e.g. the scalar  $r$  in equation B3).

$$\text{Here } \underline{V}_1 = \begin{bmatrix} e^{j-3\pi(.4)} \\ e^{j-\pi(.4)} \\ e^{j\pi(.4)} \\ e^{j3\pi(.4)} \end{bmatrix}$$

and we may choose for our complete set the following vectors

$$\underline{a}_1 = \begin{bmatrix} e^{j-3\pi(.4)} \\ e^{j-\pi(.4)} \\ e^{j\pi(.4)} \\ e^{j3\pi(.4)} \end{bmatrix} \quad \underline{a}_2 = \begin{bmatrix} -e^{j-3\pi(.4)} \\ e^{j-\pi(.4)} \\ 0 \\ 0 \end{bmatrix} \quad \underline{a}_3 = \begin{bmatrix} -e^{j-3\pi(.4)} \\ 0 \\ e^{j\pi(.4)} \\ 0 \end{bmatrix} \quad \underline{a}_4 = \begin{bmatrix} -e^{j-3\pi(.4)} \\ 0 \\ 0 \\ e^{j3\pi(.4)} \end{bmatrix} \quad (3.2.4)$$

The W matrix (equation B8) has vectors  $\underline{a}_1, \underline{W}_2, \underline{W}_3, \underline{W}_4$  as columns, where

$$\underline{W}_i = (0.11A - I) \underline{a}_i s^2 + 2A \underline{a}_i s + A (0.11A - I)^{-1} A \underline{a}_i; i=2,3,4 \quad (3.2.5)$$

The elements of this matrix are complex polynomials in s of degree two, except for the first column whose elements are just complex scalars. In this case, equation (B8) can be rewritten in terms of real and imaginary parts as follows (consider a 2x2 W matrix for simplicity):

$$\begin{bmatrix} (W_{11r} + jW_{11i}) & (W_{12r} + jW_{12i}) \\ (W_{21r} + jW_{21i}) & (W_{22r} + jW_{22i}) \end{bmatrix} \begin{bmatrix} h_{1r} + jh_{1i} \\ h_{2r} + jh_{2i} \end{bmatrix} = \begin{bmatrix} 0 + j0 \\ 0 + j0 \end{bmatrix} \quad (3.2.6)$$

This may be rearranged into the following 4x4 matrix equation

$$\begin{bmatrix} W_{11r} & -W_{11i} & W_{11r} & -W_{11i} \\ W_{11i} & W_{11r} & W_{11i} & W_{11r} \\ W_{21r} & -W_{21i} & W_{22r} & -W_{22i} \\ W_{21i} & W_{21r} & W_{22i} & W_{22r} \end{bmatrix} \begin{bmatrix} h_{1r} \\ h_{1i} \\ h_{2r} \\ h_{2i} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\text{or} \quad W \underline{H} = \underline{0} \quad (3.2.7)$$

where the new W matrix and  $\underline{H}$  vector have twice the dimension indicated by equation (B8) and are now real.



From Appendix B we know that  $h_{1r} = -1$  and  $h_{1i} = 0$ , thus the  $\underline{H}$  vector is not null, and hence the determinant of the W matrix must vanish. Setting the determinant of this W matrix equal to zero results in a polynomial of twelfth degree in s being equal to zero. Now we can theoretically proceed as before. However the numerical computation of the twelfth degree polynomial coefficients is exceedingly time consuming. We will now demonstrate that instead of having to form and solve for the roots of a twelfth degree polynomial, we can instead transform the problem into one of finding the eigenvalues of a 16 x 16 matrix, which is far easier to do numerically.

We may rewrite equation (3.2.7) in the form

$$(A_1 s^2 + A_2 s + A_3) \underline{h} = \underline{0} \quad (3.2.8)$$

where  $A_1$  and  $A_2$  are 8x8 singular matrices (their first two columns are zero), and  $A_3$  is an 8x8 invertible matrix, when we consider the four element array of example two. The problem is to find the twelve values of s for which (3.2.8) holds. Letting  $y = \frac{1}{s}$  and multiplying by  $A_3^{-1}$  gives

$$(y^2 + A_1^{-1} A_2 y + A_1^{-1} A_3) \underline{h} = \underline{0} \quad (3.2.9)$$

In terms of the two state variables

$$\underline{x}_1 = \underline{h} \quad (3.2.10a)$$

$$\underline{x}_2 = y \underline{h} \quad (3.2.10b)$$

equation (3.2.9) transforms into two first order (in y) equations

$$y \underline{x}_1 = \underline{I} \underline{x}_2 \quad (3.2.11a)$$

$$y \underline{x}_2 = y^2 \underline{h} = -A_1^{-1} A_3 \underline{x}_1 - A_1^{-1} A_2 \underline{x}_2 \quad (3.2.11b)$$

Letting  $\underline{x} = \begin{bmatrix} \underline{x}_1 \\ \underline{x}_2 \end{bmatrix}$  gives

$$y \underline{x} = \begin{bmatrix} 0 & I \\ -A_1^{-1} A_3 & -A_1^{-1} A_2 \end{bmatrix} \underline{x} \quad (3.2.12)$$

Define the 16x16 matrix G by  $G = \begin{bmatrix} 0 & I \\ -A_1^{-1} A_3 & -A_1^{-1} A_2 \end{bmatrix}$

$$y \underline{x} = G \underline{x} \quad (3.2.13)$$

Thus if s satisfies equation (3.2.7),  $y = \frac{1}{s}$  will satisfy equation (3.2.13) or

$$(G - y I) \underline{x} = \underline{0} \quad (3.2.14)$$

Therefore, instead of solving for those values of s for which equation (3.2.7) holds, we may solve for the eigenvalues  $y = \frac{1}{s}$  of the matrix G. This is much simpler.

Using this approach, we found numerically

| $y = \frac{1}{s}$ | Q      | SNR   |
|-------------------|--------|-------|
| -0.0457           | 0.0644 | ----  |
| -0.0457           | 0.0644 | ----  |
| -0.0463           | 0.0636 | ----  |
| -0.0464           | 0.0642 | ----  |
| -0.0461           | 0.0638 | ----  |
| -0.0973           | 0.110  | 0.438 |
| -0.0077           | 0.110  | 0.009 |

The remaining solutions were complex.

The best SNR we can get when the Q factor is constrained to 0.11, is SNR = 0.438. For this value of SNR, the complex vector  $\underline{I}$  is given by:

$$\begin{bmatrix} -0.096 + j 0.059 \\ 0.037 - j 0.100 \\ 0.037 + j 0.100 \\ -0.096 - j 0.059 \end{bmatrix}$$

Thus we have developed a very fast numerical technique to solve for the maximum SNR an array processor can achieve subject to a constraint on the super-gain ratio. Our next major problem is to develop an adaptive algorithm which will automatically adjust the tap weights of our array processor in such a way as to maximize the SNR subject to a constraint on the super-gain ratio. For the special cases where we have a linear array of four isotropic detectors spaced  $d = 0.8\lambda$  ( $d = 0.4\lambda$ ) apart, embedded in a uniform noise field, with the signal impinging from broadside (endfire), and with Q constrained to be equal to or less than 0.08 (0.11), we expect our adaptive array processor, in the steady state to have an output SNR which is equal to (or very close to) 0.187 (0.438). We will begin considering the design of adaptive algorithms in the next chapter.

## Appendix A      Super-Gain Ratio

It is well known that for any given aperture with a sufficiently large number of degrees of freedom (e.g. for any given detector array aperture with a sufficiently large number of array elements in it), it is possible, in theory, to obtain very high gain by using those excitations which maximize the array signal-to-noise ratio (SNR) or some similar quantity. However, this high gain is obtained at the expense of having a very large super-gain ratio (i.e. the sensitivity of the array power pattern, or gain, or SNR to small variations in the array excitations and element positions is very high). In practice therefore, since the excitations and element positions can only be controlled to within certain tolerances, it is almost impossible to actually construct super-gain arrays. To find out how well we can do in practice, we should use those excitations which are derived by maximizing the array SNR subject to a constraint on the super-gain ratio.

In this derivation of the super-gain ratio, taken from Gilbert and Morgan,<sup>(20)</sup> we will let the positions of the array elements and the element excitations vary randomly about their nominal values, with the restriction that the position random displacements have a spherically symmetrical probability distribution. It will then be shown that the expected value of the power pattern equals the nominal power pattern plus a background power level. The ratio of background power level to the nominal power pattern is directly proportional to the super-gain ratio.

### Statistical Formulation of the Super-Gain Ratio

Consider an antenna array of  $N$  elements. Each element has the same directivity pattern  $\underline{s}(\underline{r}_o)$ , where  $\underline{r}_o$  is a unit vector representing some spatial direction, and  $\underline{s}(\underline{r}_o)$  is a complex-valued vector function giving the amplitude, phase, and polarization of the radiation field over a large sphere centered at the element. For acoustic fields,  $\underline{s}(\underline{r}_o)$  is a scalar function.

The overall array directivity pattern is given by

$$\underline{p}(\underline{r}_o) = \underline{s}(\underline{r}_o) \sum_{k=1}^N J_k e^{+j \underline{k} \cdot \underline{R}_k} \quad (A 1)$$

where  $J_k$  is the complex excitation (amplitude and phase),  $k$  is the wave-number, and  $\underline{R}_k$  is a position vector from the origin to the location of the  $k^{\text{th}}$  element in the array. As usual for arrays, the pattern may be split into the element directivity pattern times the array factor  $f(\underline{r}_o)$  where

$$f(\underline{r}_o) = \sum_{k=1}^N J_k e^{+j k \underline{R}_k \cdot \underline{r}_o} \quad (\text{A } 2)$$

Note that the electric field  $\underline{E}(\underline{r}_o)$  is proportional to the array directivity pattern, i. e. the electric field strength at a point  $R \underline{r}_o$  is, for large  $R$ , proportional to

$$\frac{\underline{s}(\underline{r}_o) f(\underline{r}_o)}{R}$$

Consequently the radiated power is proportional to

$$|\underline{s}(\underline{r}_o)|^2 |f(\underline{r}_o)|^2$$

The power directivity pattern is defined as

$$\Phi(\underline{r}_o) \equiv |\underline{s}(\underline{r}_o)|^2 |f(\underline{r}_o)|^2 \quad (\text{A } 3)$$

Note that for isotropic radiators  $\underline{s}(\underline{r}_o) \equiv 1$ .

We will now assume that the excitation coefficients and the positions of the elements have some random variations about their mean or nominal values. Let

$$J_k = I_k + a_k \quad (\text{A } 4)$$

$$\underline{R}_k = \underline{r}_k + \underline{\rho}_k \quad (\text{A } 5)$$

where  $I_k$  is the nominal value of the excitation current, the  $a_k$ 's are independent random complex variables with zero mean,  $\underline{r}_k$  is the nominal value of the position vector, the  $\underline{\rho}_k$ 's are independent random vectors

with mean  $\{0, 0, 0\}$ , and all the  $\underline{\rho}_k$ 's have the same statistical distribution.

We can now find the expected values of the field and power patterns as follows:

$$\begin{aligned} E \{ \underline{s}(\underline{r}_0) f(\underline{r}_0) \} &= \underline{s}(\underline{r}_0) \sum_{k=1}^N E \{ I_k + a_k \} E \left\{ e^{jk \underline{r}_k \cdot \underline{r}_0} e^{jk \underline{\rho}_k \cdot \underline{r}_0} \right\} \\ &= \underline{s}(\underline{r}_0) \sum_{k=1}^N I_k e^{jk \underline{r}_k \cdot \underline{r}_0} E \left\{ e^{jk \underline{\rho}_k \cdot \underline{r}_0} \right\} \\ &= E \left\{ e^{jk \underline{\rho} \cdot \underline{r}_0} \right\} \underline{s}(\underline{r}_0) f_o(\underline{r}_0) \end{aligned} \quad (A 6)$$

where  $\underline{\rho}$  is a random vector having the same distribution as the  $\underline{\rho}_k$ 's, and  $f_o(\underline{r}_0)$  is the nominal array factor which results when the excitation coefficients and positions equal their nominal values.

The norm of the array factor may be written

$$\begin{aligned} |f(\underline{r}_0)|^2 &= \sum_{k=1}^N J_k e^{jk \underline{R}_k \cdot \underline{r}_0} \sum_{\ell=1}^N J_\ell^* e^{-jk \underline{R}_\ell \cdot \underline{r}_0} \\ &= \sum_{k=1}^N \sum_{\ell=1}^N (I_k + a_k) (I_\ell^* + a_\ell^*) e^{jk(\underline{r}_k + \underline{\rho}_k) \cdot \underline{r}_0} e^{-jk(\underline{r}_\ell + \underline{\rho}_\ell) \cdot \underline{r}_0} \\ &\quad k \neq \ell \\ &\quad + \sum_{k=1}^N (I_k + a_k) (I_k^* + a_k^*) \end{aligned}$$

Taking expected values and recalling that the random variables are independent

$$E \left\{ |f(\underline{r}_0)|^2 \right\} = \sum_{k=1}^N \sum_{\ell=1}^N I_k I_\ell^* e^{jk(\underline{r}_k - \underline{r}_\ell) \cdot \underline{r}_0} \left| E \left\{ e^{jk \underline{\rho} \cdot \underline{r}_0} \right\} \right|^2$$

$$+ \sum_{k=1}^N |I_k|^2 + \sum_{k=1}^N E \left\{ |a_k|^2 \right\}$$

If we now add and subtract the terms with  $k=l$  which were left out of the double sum we get

$$E \left\{ |f(\underline{r}_o)|^2 \right\} = E \left\{ \left| e^{jk \underline{\rho} \cdot \underline{r}_o} \right|^2 |f_o(\underline{r}_o)|^2 + \sum_{k=1}^N E \left\{ |a_k|^2 \right\} \right. \\ \left. + \left[ 1 - \left| E \left\{ e^{jk \underline{\rho} \cdot \underline{r}_o} \right\} \right|^2 \right] \sum_{k=1}^N |I_k|^2 \right\} \quad (A 7)$$

Multiplying through by the power pattern  $|s(\underline{r}_o)|^2$  of a single element gives the expression for the expected power pattern of the array, namely

$$E \left\{ \Phi(\underline{r}_o) \right\} = E \left\{ \left| e^{jk \underline{\rho} \cdot \underline{r}_o} \right|^2 \Phi_o(\underline{r}_o) + |s(\underline{r}_o)|^2 \sum_{k=1}^N E \left\{ |a_k|^2 \right\} \right. \\ \left. + \left[ 1 - \left| E \left\{ e^{jk \underline{\rho} \cdot \underline{r}_o} \right\} \right|^2 \right] |s(\underline{r}_o)|^2 \sum_{k=1}^N |I_k|^2 \right\} \quad (A 8)$$

where the power pattern of the nominal array is

$$\Phi_o(\underline{r}_o) = |s(\underline{r}_o)|^2 |f_o(\underline{r}_o)|^2$$

Note that in the special case where the positions of the elements are known exactly, implying that the vectors  $\underline{\rho}_k$  are all identically zero, the general result (A8) reduces to

$$E \left\{ \Phi(\underline{r}_o) \right\} = \Phi_o(\underline{r}_o) + |\underline{s}(\underline{r}_o)|^2 \sum_{k=1}^N E \left\{ |a_k|^2 \right\} \quad (A 9)$$

Equation (A9) has a simple physical interpretation. It asserts that the expected power pattern is the power pattern of the nominal array, plus a "background" power level which has the same dependence on direction as the pattern of an individual radiator, and is proportional to the sum of the mean-square errors of the excitation coefficients. In order to have the over-all pattern be a good approximation to the nominal pattern  $\Phi_o(\underline{r}_o)$ , it is necessary to hold the expected value of the background power well below the maximum value of  $\Phi_o(\underline{r}_o)$ .

If the displacements are not identically zero, Gilbert and Morgan

evaluate  $E \left\{ e^{j k \underline{\rho} \cdot \underline{r}_o} \right\}$  by assuming that the statistical distribution of  $\underline{\rho}$  is spherically symmetric, i. e. if we denote the spherical coordinates of  $\underline{\rho}$  by  $(\rho, \theta, \phi)$  then the joint probability distribution function  $p(\rho, \theta, \phi)$

depends only upon  $\rho$ . In this case the value of  $E \left\{ e^{j k \underline{\rho} \cdot \underline{r}_o} \right\}$  turns out to be independent of  $\underline{r}_o$ , and we can define a parameter  $\delta^2$  (independent of  $\underline{u}$ ) by

$$\delta^2 \equiv \left[ E \left\{ e^{j k \underline{\rho} \cdot \underline{r}_o} \right\} \right]^{-2} - 1 \quad (A10)$$

From equations (A8) and (A10) we obtain the expected power pattern for a spherically symmetric distribution of element displacements, namely

$$(1 + \delta^2) E \left\{ \Phi(\underline{r}_o) \right\} = \Phi_o(\underline{r}_o) + |\underline{s}(\underline{r}_o)|^2 \left[ (1 + \delta^2) \sum_{k=1}^N E \left\{ |a_k|^2 \right\} + \delta^2 \sum_{k=1}^N |I_k|^2 \right]$$

Again the expected pattern turns out to be the nominal pattern plus a background level with the same distribution as the pattern of a single element.

The problem is next idealized somewhat by assuming that the excitation coefficients  $J_k$  can all be controlled to the same relative accuracy, i. e. we suppose there exists a small number  $\epsilon$  such that



$$E \left\{ |a_k|^2 \right\} = \epsilon^2 |I_k|^2, \quad k = 1, 2, \dots, N. \quad (A12)$$

Then (B11) becomes

$$(1+\delta^2) E \left\{ \Phi(\underline{r}_o) \right\} = \Phi_o(\underline{r}_o) + |\underline{s}(\underline{r}_o)|^2 \left[ (1+\delta^2) \epsilon^2 + \delta^2 \right] \sum_{k=1}^N |I_k|^2 \quad (A13)$$

This expression includes the effects of both excitation and position errors.

If we define  $\Delta^2 \equiv [ (1+\delta^2) \epsilon^2 + \delta^2 ]$ , then the ratio of background power level to the average nominal power level is

$$\frac{\Delta^2 |\underline{s}(\underline{r}_o)|^2 \sum_{k=1}^N |I_k|^2}{\int_{\Omega} I_o(\underline{r}_o) d\Omega} = \frac{\Delta^2 |\underline{s}(\underline{r}_o)|^2 \sum_{k=1}^N |I_k|^2}{\int_{\Omega} |\underline{s}(\underline{r}_o)|^2 \left| \sum_{k=1}^N I_k e^{j\mathbf{k} \cdot \underline{r}_k \cdot \underline{r}_o} \right|^2 d\Omega} \quad (A14)$$

For isotropic radiators  $|\underline{s}(\underline{r}_o)|^2 = 1$ , so that the ratio becomes

$$\frac{\Delta^2 \sum_{k=1}^N |I_k|^2}{\int_{\Omega} \left| \sum_{k=1}^N I_k e^{j\mathbf{k} \cdot \underline{r}_k \cdot \underline{r}_o} \right|^2 d\Omega} = \Delta^2 Q$$

where

$$Q \equiv \frac{\sum_{k=1}^N |I_k|^2}{\int_{\Omega} \left| \sum_{k=1}^N I_k e^{j\mathbf{k} \cdot \underline{r}_k \cdot \underline{r}_o} \right|^2 d\Omega}$$

Using the vector notation of section 2.1 (see equations (2.1.1) and (2.1.4)) we may rewrite Q as

$$Q = \frac{\underline{I}^* \underline{I}}{\int_{\Omega} |\underline{I}^* V|^2 d\Omega} \quad (A15)$$

Q is a positive real number, known as the super-gain ratio, and is a measure of the sensitivity of the pattern to random errors in the excitations and positions of the array elements. Since in practice  $\Delta^2$  is never zero, an array with too large a value of Q is unacceptable.

Although Q has been derived as a result of statistical considerations, it can also be interpreted in terms of the efficiency of the array as an energy radiator. If we imagine the array elements to have a certain ohmic resistance, and the excitation coefficients to correspond to the element currents, then  $\underline{I}^* \underline{I}$  is a measure of the power which is lost in the form of heat, and Q is the ratio of dissipated power to average nominal power. Thus a large value of Q corresponds to high ohmic losses for a given amount of radiated power.

## Appendix B Maximization of SNR Subject to a Constraint

We will find the value of  $\underline{x}$  that maximizes  $\frac{\underline{x}^* \underline{C} \underline{x}}{\underline{x}^* \underline{A} \underline{x}}$  subject to the constraint  $\frac{\underline{x}^* \underline{x}}{\underline{x}^* \underline{B} \underline{x}} = q = \text{a real constant}$ , where  $\underline{A}$ ,  $\underline{B}$ , and  $\underline{C}$  are Hermitian positive definite matrices, and  $\underline{C} \equiv \underline{a}_1 \underline{a}_1^*$ . This appendix represents work done by Lo, Lee and Lee (19).

Introducing a real scalar Lagrange multiplier  $\lambda$ , the solution can be obtained by differentiating  $L$  with respect to  $\underline{x}$ , and setting the result equal to zero, where

$$L = \frac{\underline{x}^* \underline{C} \underline{x}}{\underline{x}^* \underline{A} \underline{x}} + \lambda \frac{\underline{x}^* \underline{x}}{\underline{x}^* \underline{B} \underline{x}} \quad (\text{B } 1)$$

Thus

$$\delta \underline{x}^* \left\{ \frac{\underline{C} \underline{x} (\underline{x}^* \underline{A} \underline{x}) - \underline{A} \underline{x} (\underline{x}^* \underline{C} \underline{x})}{(\underline{x}^* \underline{A} \underline{x})^2} + \frac{\underline{x} (\underline{x}^* \underline{B} \underline{x}) \lambda - \underline{B} \underline{x} (\underline{x}^* \underline{x}) \lambda}{(\underline{x}^* \underline{B} \underline{x})^2} \right\} \\ + \left\{ \frac{(\underline{x}^* \underline{A} \underline{x}) \underline{x}^* \underline{C} - (\underline{x}^* \underline{C} \underline{x}) \underline{x}^* \underline{A}}{(\underline{x}^* \underline{A} \underline{x})^2} + \frac{\lambda (\underline{x}^* \underline{B} \underline{x}) \underline{x}^* - \lambda (\underline{x}^* \underline{x}) \underline{x}^* \underline{B}}{(\underline{x}^* \underline{B} \underline{x})^2} \right\} \delta \underline{x}$$

= 0

Since  $\underline{A}$ ,  $\underline{B}$ , and  $\underline{C}$  are Hermitian

$$(\underline{x}^* \underline{A} \delta \underline{x}) = (\delta \underline{x}^* \underline{A} \underline{x})^*$$

$$(\underline{x}^* \underline{B} \delta \underline{x}) = (\delta \underline{x}^* \underline{B} \underline{x})^*$$

$$(\underline{x}^* \underline{C} \delta \underline{x}) = (\delta \underline{x}^* \underline{C} \underline{x})^*$$

Making this substitution in the second term of the last equation results in the second term becoming

$$\left[ \delta \underline{x}^* \left\{ \frac{C \underline{x} (\underline{x}^* A \underline{x}) - A \underline{x} (\underline{x}^* C \underline{x})}{(\underline{x}^* A \underline{x})^2} + \frac{\underline{x} (\underline{x}^* B \underline{x}) \lambda - B \underline{x} (\underline{x}^* \underline{x}) \lambda}{(\underline{x}^* B \underline{x})^2} \right\} \right]^*$$

Note that the terms inside the braces are equal to the terms inside the braces in the first term of the last equation. Thus, the overall equation is of the form

$$\delta \underline{x}^* \underline{y} + (\delta \underline{x}^* \underline{y})^* = 0$$

Since this equation must be true for all possible values of the real and imaginary parts of  $\delta \underline{x}$ , this implies  $\underline{y} = \underline{0}$ .

Thus

$$\frac{C \underline{x} (\underline{x}^* A \underline{x}) - A \underline{x} (\underline{x}^* C \underline{x})}{(\underline{x}^* A \underline{x})^2} + \frac{\underline{x} (\underline{x}^* B \underline{x}) \lambda - B \underline{x} (\underline{x}^* \underline{x}) \lambda}{(\underline{x}^* B \underline{x})^2} = \underline{0} \quad (B2)$$

But  $C \equiv a_1 a_1^*$  and we can assume  $\underline{x}$  is normalized to 1, i.e.  $\underline{x}^* \underline{x} = 1$

because both the function we are maximizing  $\frac{\underline{x}^* C \underline{x}}{\underline{x}^* A \underline{x}}$  and the constraint

$\frac{\underline{x}^* \underline{x}}{\underline{x}^* B \underline{x}}$  are independent of the magnitude of  $\underline{x}$ . Multiplying equation (B2)

by  $(\underline{x}^* A \underline{x})$ , letting  $C = a_1 a_1^*$  in the first term, and multiplying the third and fourth terms by  $\underline{x}^* \underline{x} = 1$ , gives

$$a_1 a_1^* \underline{x} - \frac{A \underline{x} (\underline{x}^* C \underline{x})}{(\underline{x}^* A \underline{x})} + \frac{\lambda \underline{x} (\underline{x}^* \underline{x}) (\underline{x}^* A \underline{x})}{(\underline{x}^* B \underline{x})} - \frac{\lambda B \underline{x} (\underline{x}^* \underline{x})^2 (\underline{x}^* A \underline{x})}{(\underline{x}^* B \underline{x})^2} = \underline{0}$$

since  $q = \frac{\underline{x}^* \underline{x}}{\underline{x}^* B \underline{x}}$  we have

$$\underline{a}_1 (\underline{a}_1^* \underline{x}) - \frac{A \underline{x} (\underline{x}^* C \underline{x})}{(\underline{x}^* A \underline{x})} + \lambda q (\underline{x}^* A \underline{x}) \underline{x} - \lambda q^2 (\underline{x}^* A \underline{x}) B \underline{x} = 0$$

Combining terms

$$(\underline{a}_1^* \underline{x}) \underline{a}_1 = \left[ \frac{(\underline{x}^* C \underline{x})}{(\underline{x}^* A \underline{x})} A - \lambda q (\underline{x}^* A \underline{x}) I + \lambda q^2 (\underline{x}^* A \underline{x}) B \right] \underline{x}$$

Multiplying by the real scalar  $\frac{(\underline{x}^* A \underline{x})}{(\underline{x}^* C \underline{x})}$  gives

$$\frac{(\underline{a}_1^* \underline{x}) (\underline{x}^* A \underline{x})}{(\underline{x}^* C \underline{x})} \underline{a}_1 = \left[ A - \lambda q \frac{(\underline{x}^* A \underline{x})^2 I}{(\underline{x}^* C \underline{x})} + \lambda q^2 \frac{(\underline{x}^* A \underline{x})^2}{(\underline{x}^* C \underline{x})} B \right] \underline{x}$$

$$\text{Define } r \equiv \frac{(\underline{a}_1^* \underline{x}) (\underline{x}^* A \underline{x})}{(\underline{x}^* C \underline{x})} \quad \text{a complex scalar} \quad (B 3)$$

$$s = \frac{\lambda q (\underline{x}^* A \underline{x})^2}{(\underline{x}^* C \underline{x})} \quad \text{a real scalar} \quad (B 4)$$

$$\text{thus } r \underline{a}_1 = [A - sI + q s B] \underline{x}$$

The solution for  $\underline{x}$  is

$$\underline{x} = r K^{-1} \underline{a}_1 \quad (B 5)$$

where  $r$  is a complex scalar, depending upon  $\underline{x}$ , and  $K \equiv [A - sI + q s B]$  is a Hermitian matrix which also depends upon  $\underline{x}$ .

In addition to equation (B 5), the constraint equation must also be satisfied, thus

$$q = \frac{\underline{x}^* \underline{x}}{\underline{x}^* B \underline{x}} \quad (B 6)$$

Since only the direction of  $\underline{x}$  and not its magnitude (we showed its magnitude could be assumed equal to unity) is of interest, the scalar  $r$  which multiplies all components of  $\underline{x}$  may be disregarded. The only unknown, then, in the simultaneous solution of equations (B5) and (B6) is the real scalar  $s$ , which is proportional to the Lagrange multiplier  $\lambda$ . Inserting (B5) into (B6) one obtains a characteristic equation for  $s$ .

$$q = \frac{\underline{a}_1^* K^{-1} K^{-1} \underline{a}_1}{\underline{a}_1^* K^{-1} B K^{-1} \underline{a}_1}$$

this may be rewritten in the form

$$\begin{aligned} \underline{a}_1^* K^{-1} q B K^{-1} \underline{a}_1 - \underline{a}_1^* K^{-1} I K^{-1} \underline{a}_1 &= 0 \\ \underline{a}_1^* K^{-1} [q B - I] K^{-1} \underline{a}_1 &= 0 \end{aligned} \quad (B 7)$$

Because the unknown  $s$  is contained in  $K$ , a direct numerical solution of (B7) is very difficult. However, Lo, Lee and Lee observed that equation (B7) states that the vector  $\underline{a}_1$  is orthogonal to the vector  $K^{-1} [q B - I] K^{-1} \underline{a}_1$ . Thus the vector  $K^{-1} [q B - I] K^{-1} \underline{a}_1$  must lie in the space orthogonal to the space spanned by  $\underline{a}_1$ . A complete set  $\{a_n\}$  with  $\underline{a}_1$  as one of its elements can be easily constructed, e.g. if

$$\underline{a}_1 = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \quad \text{we may choose}$$

$$\underline{a}_2 = \begin{bmatrix} -1 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \underline{a}_3 = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \dots \quad \underline{a}_n = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}$$

The vector  $K^{-1} [qB - I] K^{-1} \underline{a}_1$  must be a linear combination of the vectors  $\underline{a}_2, \underline{a}_3, \dots, \underline{a}_N$ . Let it be

$$K^{-1} [qB - I] K^{-1} \underline{a}_1 = \sum_{n=2}^N h_n \underline{a}_n$$

which yields

$$\underline{a}_1 = \sum_{n=2}^N h_n K [qB - I]^{-1} K \underline{a}_n$$

rearranging gives

$$\sum_{n=2}^N [A + s(qB - I)] [qB - I]^{-1} [A + s(qB - I)] \underline{a}_n h_n - \underline{a}_1 = \underline{0}$$

$$\sum_{n=2}^N \left\{ A(qB - I)^{-1} A \underline{a}_n + A s \underline{a}_n + s A \underline{a}_n + s^2 (qB - I) \underline{a}_n \right\} h_n - \underline{a}_1 = \underline{0}$$

$$\underline{a}_1(-1) + \sum_{n=2}^N \left\{ s^2 (qB - I) \underline{a}_n + 2s A \underline{a}_n + A(qB - I)^{-1} A \underline{a}_n \right\} h_n = \underline{0}$$

$$\text{or} \quad -1 \underline{a}_1 + h_2 \underline{W}_2 + h_3 \underline{W}_3 + \dots + h_N \underline{W}_N = \underline{0}$$

in matrix form

$$W \underline{H} = \underline{0} \tag{B8}$$

where  $W$  is a matrix with in general, complex vectors  $\underline{a}_1, \underline{W}_2, \underline{W}_3, \dots, \underline{W}_N$  as columns, i.e.  $\underline{W}_n = s^2 (qB - I) \underline{a}_n + 2s A \underline{a}_n + A(qB - I)^{-1} A \underline{a}_n$   
 $n=2, 3, \dots, N$

$$\underline{W}_1 = \underline{a}_1$$

and

$$\underline{H} = \begin{bmatrix} -1 \\ h_2 \\ \vdots \\ h_N \end{bmatrix}$$

Since  $\underline{H}$  is not a null vector, the determinant of  $\underline{W}$  in equation (B8) must vanish, i.e.

$$\det [\underline{a}_1, \underline{W}_2, \dots, \underline{W}_N] = 0 \quad (B 9)$$

This results in a (sometimes complex) polynomial of degree  $2(N-1)$  in the unknown  $s$ , and thus the roots can be numerically determined. One

of them will give the absolute maximum of  $\frac{\underline{x}^* C \underline{x}}{\underline{x}^* A \underline{x}}$ , because once the

possible value of  $s$  have been found, the direction of  $\underline{x}$  can be found from equation (B 5) and the problem is solved.



## CHAPTER 4

### Minimization of the Mean-Squared-Error (MSE)

#### Subject to One Linear Constraint

Our objective is to consider an adaptive algorithm which will maximize the SNR subject to a constraint on the super-gain ratio when unknown interfering noise is present. Because the SNR and super-gain ratio are nonlinear quantities, it is difficult to prove convergence of our algorithm to the optimal solution, or to analytically find the algorithm's rate of convergence. Thus, for the purpose of mathematical tractability (the nonlinear algorithm will be simulated on a computer to obtain some numerical indication of convergence and convergence rate in chapter six), and because (1) the criterion of minimizing the MSE is important in its own right (2) linear constraints may appear in similar problems (3) nonlinear constraints are approximately linear near the solution point and (4) the projection method used in the linear case is also applicable to the nonlinear case, we will consider in this chapter an adaptive algorithm which minimizes the MSE subject to a linear constraint. Specifically, we will find the Lagrange solution to the problem of minimizing the MSE subject to a linear constraint and then prove that an algorithm of the form  $\underline{W}_{j+1} = \underline{W}_j - k P \nabla_{\underline{W}_j} (\text{MSE})$  converges to the Lagrange solution, when the gradient  $\nabla_{\underline{W}_j} (\text{MSE})$  is (1) known exactly, (2) estimated, and (3) estimated by an estimate which contains additive noise.

# Section 4.1 Derivation of Mean Squared Error and Constraint Equation

The processor configuration is shown in Fig 4.1.1 where  $\Delta$  represents a time delay,  $\underline{s}_j \equiv \text{col}(s_{1j}, s_{2j}, \dots, s_{nj})$  is the stochastic signal at the outputs of the tapped delay lines at time (iteration)  $j$ , the  $W$ 's are the multiplicative tap weights, and  $d_j$  is some known scalar function of the vector  $\underline{s}_j$ , i.e.  $d_j$  represents the desired array output at time  $j$ .

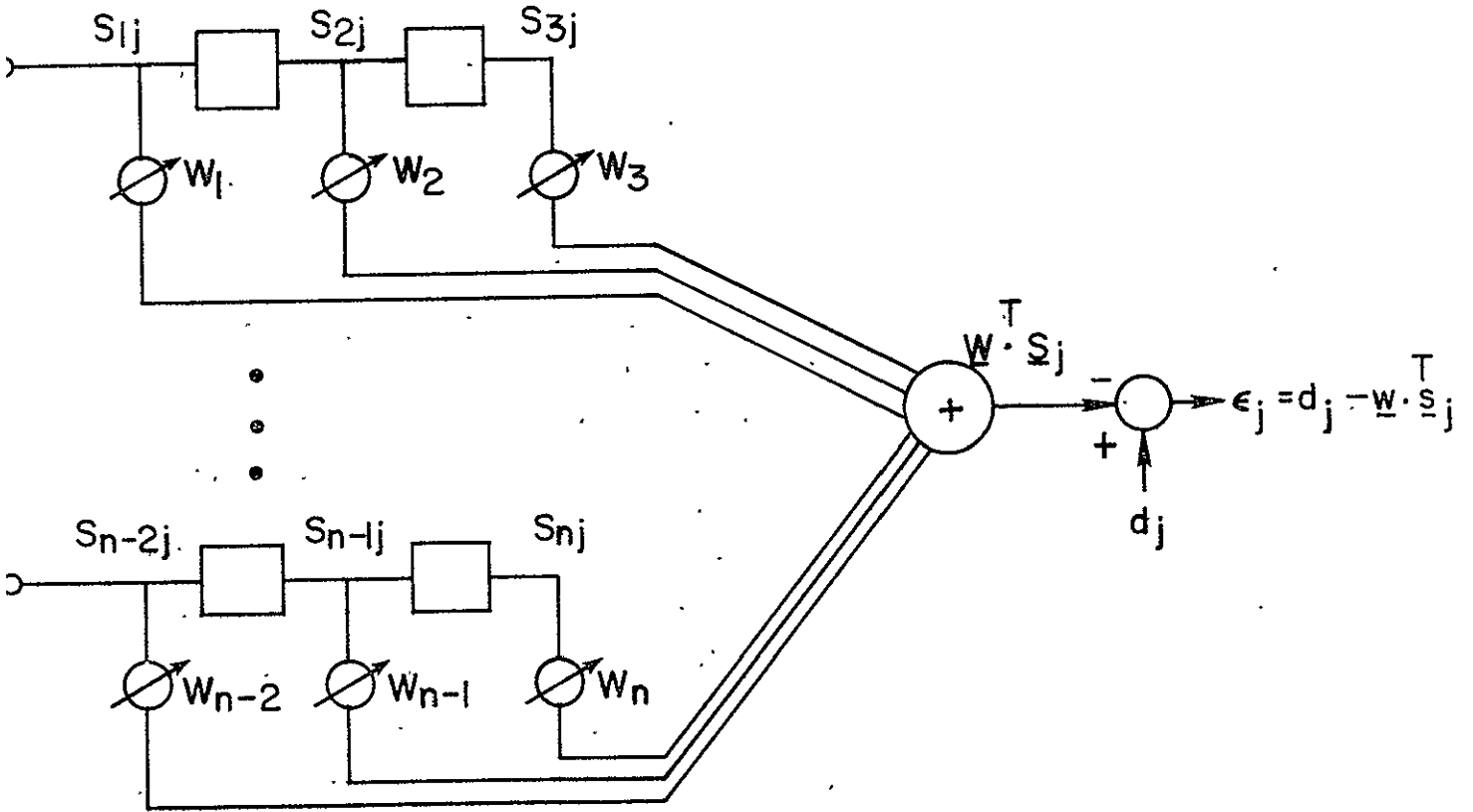


Fig. 4.1.1 Processor Configuration

From Fig 4.1.1 we have

$$\epsilon_j = d_j - \underline{W}^T \underline{s}_j \quad (4.1.1)$$

$$\epsilon_j^2 = d_j^2 - 2d_j \underline{s}_j^T \underline{W} + \underline{W}^T \underline{s}_j \underline{s}_j^T \underline{W} \quad (4.1.2)$$

When the input signal can be regarded as a stationary, ergodic random process, then

$$E\{\underline{s}_j\} = \underline{s} \text{ and } E\{d_j\} = d$$

Our problem is to devise an algorithm that will adjust the weights to their LMS value subject to a linear constraint. Toward this end we have already found an expression (equation 4.1.2) for the MSE, and the remainder of this section will be devoted to finding expressions for the minimum value of the MSE when we have no constraint, mention of an adaptive algorithm that will automatically adjust the tap weights to their unconstrained LMS values, and writing an expression for any arbitrary linear constraint on  $\underline{W}$ . Taking the expected value of equation (4.1.2) gives

$$E\{\epsilon_j^2\} \equiv \overline{\epsilon_j^2} = \overline{d_j^2} - 2 \underline{\phi}^T(\underline{s}, d) \underline{W}_j + \underline{W}_j^T \underline{\phi}(\underline{s}, \underline{s}) \underline{W}_j \quad (4.1.3)$$

where

$$\underline{\phi}(\underline{s}, d) \equiv E\{\underline{s}_j d_j\} = \begin{bmatrix} E\{s_{1j} d_j\} \\ \vdots \\ E\{s_{nj} d_j\} \end{bmatrix} \quad (4.1.4)$$

$$\underline{\phi}(\underline{s}, \underline{s}) \equiv E\{\underline{s}_j \underline{s}_j^T\} \quad (4.1.5)$$

Taking the gradient of  $\overline{\epsilon_j^2}$  yields

$$\nabla^T(\overline{\epsilon_j^2}) = -2 \underline{\phi}^T(\underline{s}, d) + 2 \underline{W}^T \underline{\phi}(\underline{s}, \underline{s}) \quad (4.1.6)$$

To find the least-mean-square (LMS) set of weights,  $\underline{W}_{LMS}$ , that minimizes  $\overline{\epsilon_j^2}$  when there is no constraint, we set  $\nabla(\overline{\epsilon_j^2}) = 0$ . Thus

$$\underline{\phi}^T(\underline{s}, d) = \underline{W}_{LMS}^T \underline{\phi}(\underline{s}, \underline{s}) \quad (4.1.7a)$$

$$\underline{W}_{LMS}^T = \underline{\phi}^T(\underline{s}, d) \underline{\phi}^{-1}(\underline{s}, \underline{s}) \quad (4.1.7b)$$

The LMS error is achieved by choosing the optimal weight vector given by equation (4.1.7b). An expression for the minimum mean-square error may be obtained by substituting (4.1.7a) into (4.1.3)

$$\overline{\epsilon_{\min}^2} \equiv \min(\overline{\epsilon_j^2}) = \overline{d_j^2} - \underline{W}_{LMS}^T \underline{\phi}(\underline{s}, \underline{s}) \underline{W}_{LMS} \quad (4.1.8)$$

Note that  $\min(\overline{\epsilon_j^2})$  is independent of  $j$  ( $\overline{d_j^2}$  is independent of  $j$ ).

Widrow, Lucky and others (12)-(18) have investigated adaptive algorithms which automatically adjust the tap weights to their unconstrained LMS values. One such algorithm is given by

$$\underline{W}_{j+1} = \underline{W}_j - k \nabla(\overline{\epsilon_j^2}) \quad (4.1.9)$$

Substituting (4.1.6) into (4.1.9) gives

$$\underline{W}_{j+1} = \underline{W}_j + 2k \underline{\phi}(\underline{s}, d) - 2k \underline{\phi}(\underline{s}, \underline{s}) \underline{W}_j \quad (4.1.10)$$

Note that equation (4.1.10) is a linear equation in  $\underline{W}$ . This means we can easily solve for  $\lim_{j \rightarrow \infty} \underline{W}_j$  and other quantities of interest, and it is the main

reason we are using minimum mean-square error as our criterion. The abovementioned researchers have proven that by using the algorithm of equation (4.1.9)  $\underline{W}_j$  converges to  $\underline{W}_{LMS}$ .

Any arbitrary linear constraint on  $\underline{W}$  can be written in the form

$$\underline{W}^T \underline{n}_1 - a \geq 0 \quad (4.1.11)$$

where  $\underline{n}_1$  is a unit normal to the hyperplane  $\underline{W}^T \underline{n}_1 - a = 0$ .

Our problem now is to (1) find the optimum value of the weights,  $\underline{W}_{opt}$ , which yields the minimum MSE (equation 4.1.2) subject to the constraint (4.1.11) and (2) devise an adaptive algorithm, similar to (4.1.9) which will make the tap weights  $\underline{W}$  converge to this  $\underline{W}_{opt}$ . The next section attacks the first problem.

## Section 4.2 Analytic (Lagrange) Solution

In this section we will use a Lagrange multiplier technique to find the optimum value of the weights  $\underline{W}_{opt}$ , which yields the minimum mean-square-error subject to the linear constraint (4.1.11).

Let us first rewrite equation (4.1.3) for  $\overline{\epsilon_j^2}$  as follows.

Substituting (4.1.7a) and (4.1.8) into (4.1.3) gives

$$\overline{\epsilon_j^2} = \left[ \overline{\epsilon_{min}^2} + \underline{W}_{LMS}^T \phi(\underline{s}, \underline{s}) \underline{W}_{LMS} \right] - 2 \underline{W}_{LMS}^T \phi(\underline{s}, \underline{s}) \underline{W} + \underline{W}^T \phi(\underline{s}, \underline{s}) \underline{W}$$

But

$$\underline{W}_{LMS}^T \phi(\underline{s}, \underline{s}) \underline{W} = \underline{W}^T \phi(\underline{s}, \underline{s}) \underline{W}_{LMS}$$

Thus

$$\overline{\epsilon_j^2} = \overline{\epsilon_{min}^2} + (\underline{W}^T \underline{W}_{LMS}^T) \phi(\underline{s}, \underline{s}) (\underline{W} - \underline{W}_{LMS}) \quad (4.2.1)$$

The problem is to maximize (4.2.1) subject to (4.1.11). Let us investigate what the solution looks like both graphically and analytically. Graphically we have

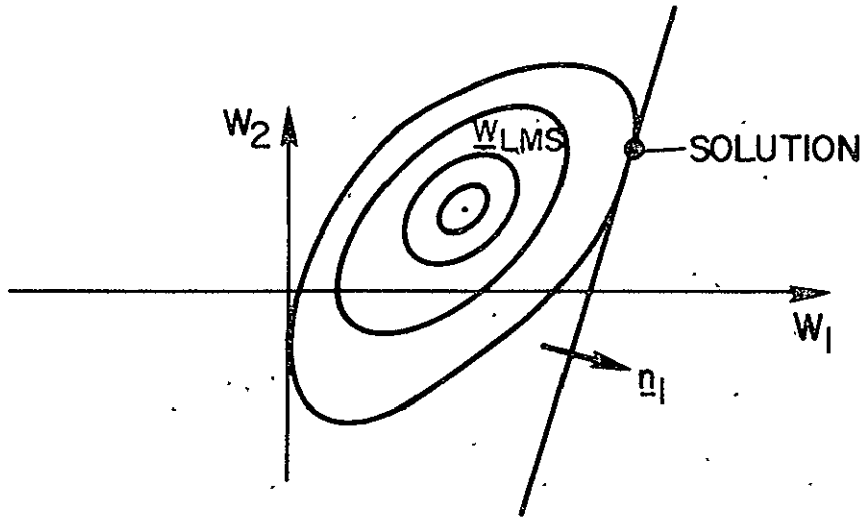


Fig. 4.2.1 Typical MSE level curves and constraint

Since the objective function is quadratic, the solution is either:

1.  $\underline{W} = \underline{W}_{LMS}$  or

2.  $\underline{W}$  = the solution to the Lagrange multiplier problem when (4.1.11) holds as an equality, i.e.  $\underline{W}^T \underline{n}_1 - a = 0$ .

We are only interested in case (2) in this section, because the algorithms of Widrow and Lucky will work in case (1).

Analytically we must minimize

$$\overline{\epsilon_j^2} = \overline{\epsilon_{min}^2} + (\underline{W}^T - \underline{W}_{LMS}^T) \phi(\underline{W} - \underline{W}_{LMS})$$

subject to the constraint

$$\underline{W}^T \underline{n}_1 - a = 0$$

The Lagrange technique yields

$$L = \overline{\epsilon_{min}^2} + (\underline{W}^T - \underline{W}_{LMS}^T) \phi(\underline{W} - \underline{W}_{LMS}) + \alpha \left[ \underline{W}^T \underline{n}_1 - a \right] \quad (4.2.2)$$

Taking differentials with respect to  $\underline{W}$  we have

$$\begin{aligned} \delta L = (\delta \underline{W}^T) \phi \underline{W} + \underline{W}^T \phi (\delta \underline{W}) - (\delta \underline{W}^T) \phi \underline{W}_{LMS} - \underline{W}_{LMS}^T \phi (\delta \underline{W}) \\ + \alpha (\delta \underline{W}^T) \underline{n}_1 = 0 \end{aligned} \quad (4.2.3)$$

But

$$\left\{ (\delta \underline{W}^T) [\phi \underline{W} - \phi \underline{W}_{LMS}] \right\}^T = [\underline{W}^T \phi - \underline{W}_{LMS}^T \phi] (\delta \underline{W})$$

(4.2.3) may be rewritten as

$$\left( \alpha \underline{n}_1^T + 2 [\underline{W}^T \phi - \underline{W}_{LMS}^T \phi] \right) (\delta \underline{W}) = 0$$

Which must be true for all  $\delta \underline{W}$ , giving

$$\alpha \underline{n}_1^T + 2 [\underline{W}^T \phi - \underline{W}_{LMS}^T \phi] = 0 \quad (4.2.4)$$

equation (4.2.4) together with the constraint equation (4.2.1) must be solved simultaneously for  $\alpha$  and  $\underline{W}$ . Doing this yields

$$\underline{W}_{\text{optimum}}^T = \frac{(\alpha - \underline{W}_{LMS}^T \underline{n}_1)}{(\underline{n}_1^T \phi^{-1} \underline{n}_1)} \quad \underline{n}_1^T \phi^{-1} + \underline{W}_{LMS}^T \quad (4.2.5)$$

This is the analytic solution for the least mean square value of the tap weights subject to an arbitrary linear constraint. In the next section we will present an adaptive algorithm, which will, in the steady state, make the tap weights converge to this optimum value we have just found in equation (4.2.5).

### Section 4.3 Use of the Projected Gradient Algorithm to Adaptively Adjust the Tap Weights

The projected gradient algorithm that we will use is a modified version of Rosen's algorithm which is discussed briefly in Appendix B. It is advisable to read Appendix B before the following sections. The algorithm we will use to minimize the MSE subject to a linear constraint may be thought of intuitively as follows: We want to converge to the vector  $\underline{W}_{opt}$  which minimizes the MSE, which is a function of  $\underline{W}$ , subject to a linear constraint on the vector  $\underline{W}$ . Looking at Fig 4.3.1 we see intuitively that

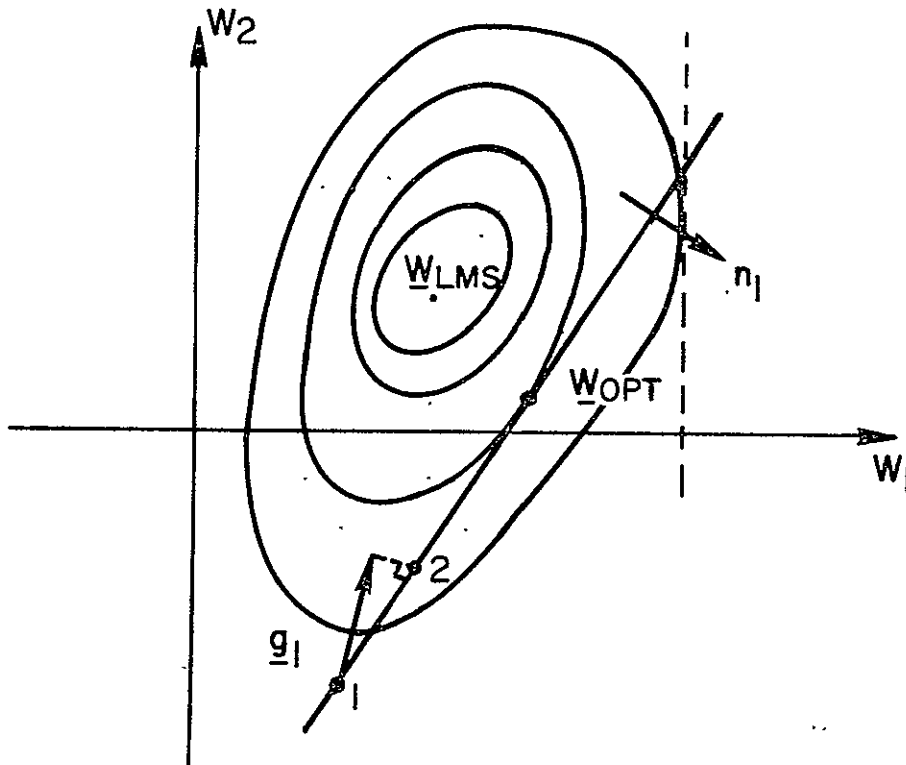


Fig. 4.3.1 Intuitive idea behind projected gradient algorithm

we can start at a point which satisfies the linear constraint, denote it by point one; find the gradient of the MSE with respect to  $\underline{W}$  at point one and "project" this gradient vector, which lies in an  $n$  dimensional vector space (in Fig 4.3.1 the  $n$  dimensional  $\underline{W}$  vector space is of dimension 2), onto the  $n-1$  (one dimensional in the diagram) dimensional subspace which is orthogonal to the one dimensional subspace spanned by the normal  $\underline{n}_1$  to the constraint surface, call this point two; and repeat the procedure



indefinitely. This procedure may converge to the constrained optimum denoted by  $\underline{W}_{opt}$  under certain conditions.

Analytically, the projected gradient algorithm is given by

$$\underline{W}_{j+1} = \underline{W}_j - k P \nabla_{\underline{W}_j} (\text{MSE})$$

where  $P$  is the projection operator  $P = I - \underline{n}_1 \underline{n}_1^T$  if we have only one constraint (see Appendix B for the more general case),  $\underline{n}_1$  is a unit vector normal to the constraint hyperplane,  $k$  is a constant which will be investigated later, and  $\nabla_{\underline{W}_j} (\text{MSE})$  is the gradient of MSE at time (iteration)  $j$ .

#### Section 4.3.1      The Algorithm, Proof of Convergence, and Bounds on the Rate of Convergence if the Gradient is Known.

Let us compute the gradient of the MSE,  $\underline{g}$ , and the gradient projection  $P \underline{g}$ . From equation (4.1.6)

$$\underline{g}^T \equiv \nabla_{\underline{W}}^T (\epsilon_j^2) = -2 \underline{\phi}^T(\underline{s}, \underline{d}) + 2 \underline{W}^T \underline{\phi}(\underline{s}, \underline{s})$$

using (4.1.7a) we get

$$\underline{g} = 2 \underline{\phi} [ \underline{W} - \underline{W}_{LMS} ] \quad (4.3.1.1)$$

The projection operator is given by

$$P = I - \underline{n}_1 \underline{n}_1^T \quad (4.3.1.2)$$

thus

$$P \underline{g} = [I - \underline{n}_1 \underline{n}_1^T] 2 \underline{\phi} [ \underline{W} - \underline{W}_{LMS} ] \quad (4.3.1.3)$$

Our algorithm is

$$\underline{W}_{j+1} = \underline{W}_j - k [I - \underline{n}_1 \underline{n}_1^T] 2 \underline{\phi} [ \underline{W}_j - \underline{W}_{LMS} ] \quad (4.3.1.4)$$

As discussed before, we will start at a point where the constraint is satisfied, and since at every iteration we are projecting  $\underline{W}$  onto a subspace where the constraint is satisfied, this implies that the constraint

equation is always, satisfied, i.e.

$$\underline{W}_j^T \cdot \underline{n}_1 = a \quad j = 0, 1, 2, \dots$$

Equations (4.3.1.4) constitute a set of  $n$  simultaneous first order difference equations. In order to solve them, we need initial conditial conditions. For our "initial" conditions, we will use the fact that the constraint must always be satisfied, and in particular must be satisfied at  $j = \infty$ , i.e.

$$\underline{W}_\infty^T \cdot \underline{n}_1 = a \quad (4.3.1.5)$$

Now equations (4.3.1.4) and (4.3.1.5) constitute a set of  $n$  first order deterministic difference equations (since  $\underline{W}$  is of dimension  $n$ ) with initial conditions. We want to investigate whether or not the sequence of  $\underline{W}$ 's converges to  $\underline{W}_{opt}$ , and if so, what is the rate of convergence?

To answer the first question, we will solve for the asymptotic value of equation (4.3.1.4)

$$\underline{W}_\infty = \underline{W}_\infty + 2k [I - \underline{n}_1 \underline{n}_1^T] \phi [\underline{W}_{LMS} - \underline{W}_\infty]$$

$$\underline{0} = [I - \underline{n}_1 \underline{n}_1^T] \phi [\underline{W}_{LMS} - \underline{W}_\infty]$$

$$\text{Let} \quad \underline{x} \equiv \underline{W}_\infty - \underline{W}_{LMS} \quad (4.3.1.6)$$

$$\text{then} \quad [I - \underline{n}_1 \underline{n}_1^T] \phi \underline{x} = \underline{0} \quad (4.3.1.7)$$

Again, since  $\underline{W}$  has  $n$  components, equations (4.3.1.7) constitute a set of  $n$  simultaneous deterministic homogeneous equations in  $n$  unknowns. The initial condition (4.3.1.5) becomes

$$\underline{n}_1^T \cdot \underline{x} = a - \underline{n}_1^T \cdot \underline{W}_{LMS} \quad (4.3.1.8)$$

Before solving (4.3.1.7) let us consider the following equations.

$$A \underline{x} = \underline{o}$$

1. A necessary and sufficient condition for the above  $n$  equations to have a nontrivial solution is that the rank of  $A$  be less than  $n$ , or equivalently, that the determinant of  $A$  be zero.

2. If the rank of  $A$  is  $r$ , where  $r < n$ , then the system of equations has exactly  $n - r$  linearly independent solutions such that every solution is a linear combination of these  $n - r$  linearly independent solutions and every linear combination of the  $n - r$  linearly independent solutions is a solution.

Let us now investigate the rank of  $[I - \underline{n}_1 \underline{n}_1^T] \phi$ . By definition, the rank of an operator is the dimension of the range space of the operator, thus

$$\text{rank} [I - \underline{n}_1 \underline{n}_1^T] = n - 1$$

For arbitrary matrices  $B$  and  $C$

$$\text{rank} (BC) \leq \min (\text{rank } B, \text{rank } C)$$

From this we may conclude that

1. Because  $\text{rank} [I - \underline{n}_1 \underline{n}_1^T] = n - 1$ , this implies there exists at least one (possibly nonunique) solution to equations (4.3.1.7).

2. If we know that the rank of  $[I - \underline{n}_1 \underline{n}_1^T] \phi$  equals  $n - 1$ , this implies there exists a unique (to within a multiplicative constant-which is unique provided the initial condition is satisfied) nontrivial solution to equations (4.3.1.7).

If  $\phi$  is invertible, then the rank of  $[I - \underline{n}_1 \underline{n}_1^T] \phi = \text{rank} [I - \underline{n}_1 \underline{n}_1^T] = n - 1$ . This follows from Halmos, <sup>(23)</sup> Theorem 3, part IV, page 92. Since  $\phi$  is a correlation matrix, it is positive semidefinite, and, in practice almost always positive definite, which implies that it is invertible. Thus equations (4.3.1.7), together with the initial conditions of equations (4.3.1.8) have a unique solution.

If  $\underline{W}_\infty = \underline{W}_{\text{optimum}}$  satisfies (4.3.1.7) and (4.3.1.8) then it is the solution. We will now verify that this is the case. From (4.2.5)

$$\underline{W}_\infty = \underline{W}_{\text{opt}} = \frac{(a - \underline{W}_{\text{LMS}}^T \underline{n}_1)}{(\underline{n}_1^T \phi^{-1} \underline{n}_1)} \phi^{-1} \underline{n}_1 + \underline{W}_{\text{LMS}}$$

$$\underline{x} = \underline{W}_\infty - \underline{W}_{\text{LMS}} = \frac{(a - \underline{W}_{\text{LMS}}^T \underline{n}_1)}{(\underline{n}_1^T \phi^{-1} \underline{n}_1)} \phi^{-1} \underline{n}_1$$

Substituting this expression for  $\underline{x}$  into (4.3.1.7) and (4.3.1.8) one sees that the equations are satisfied. Thus  $\underline{W}_\infty = \underline{W}_{\text{opt}}$  is the unique solution to equations (4.3.1.7) and (4.3.1.8).

Now that we have shown that the sequence of  $\underline{W}$ 's does converge to  $\underline{W}_{\text{opt}}$ , we will investigate the rate of convergence of the weight vectors to  $\underline{W}_{\text{opt}}$ , given by (4.2.5)

$$\underline{W}_{\text{opt}} = \frac{(a - \underline{W}_{\text{LMS}}^T \underline{n}_1)}{(\underline{n}_1^T \phi^{-1} \underline{n}_1)} \phi^{-1} \underline{n}_1 + \underline{W}_{\text{LMS}}$$

Define

$$\underline{q}_j \equiv \underline{W}_j - \underline{W}_{\text{opt}} \quad (4.3.1.9)$$

The algorithm (4.3.1.4) can be rewritten as

$$\underline{W}_{j+1} = \left[ (1-2k\phi) I + 2k \underline{n}_1 \underline{n}_1^T \phi \right] \underline{W}_j + 2k \left[ I - \underline{n}_1 \underline{n}_1^T \right] \phi \underline{W}_{\text{LMS}}$$

After some manipulation (and noting that  $[\mathbf{I} - \underline{n}_1 \underline{n}_1^T] \cdot \phi^{-1} \phi \underline{n}_1 = \underline{0}$ ) we have

$$\underline{q}_{j+1} = \left[ \mathbf{I} - 2k (\mathbf{I} - \underline{n}_1 \underline{n}_1^T) \phi \right] \underline{q}_j \quad (4.3.1.10)$$

Since  $\underline{q}_j = \underline{W}_j - \underline{W}_{\text{opt}}$ , by looking at Fig. 4.3.1 we see that  $\underline{q}_j$  always lies in the hyperplane (in the Figure this means lie along the constraint line) which is orthogonal to  $\underline{n}_1$ , hence

$$\mathbf{P} \underline{q}_j = \underline{q}_j \quad \text{for all } j \quad (4.3.1.11)$$

Thus

$$\underline{q}_{j+1} = \left[ \mathbf{I} - \underline{n}_1 \underline{n}_1^T \right] (\mathbf{I} - 2k \phi) \underline{q}_j \quad (4.3.1.12)$$

and

$$\| \underline{q}_{j+1} \| \leq \xi^{j+1} \| \underline{q}_0 \| \quad (4.3.1.13)$$

where

$$\xi \equiv \left\| \left[ \mathbf{I} - \underline{n}_1 \underline{n}_1^T \right] (\mathbf{I} - 2k \phi) \right\| \quad (4.3.1.14)$$

Let us investigate this norm. The correlation matrix  $\phi$  is a symmetric positive semidefinite, and in practice almost always positive definite, matrix with positive minimum and maximum eigenvalues  $\rho_1$  and  $\rho_N$  respectively;  $k$  is chosen to be a positive number; and  $(\mathbf{I} - \underline{n}_1 \underline{n}_1^T)$  is a projection operator as discussed previously. To bound the norm, we have

$$\xi \leq \left\| \mathbf{I} - \underline{n}_1 \underline{n}_1^T \right\| \left\| \mathbf{I} - 2k \phi \right\| \quad (4.3.1.15)$$

Since  $\mathbf{I} - \underline{n}_1 \underline{n}_1^T$  is a projection operator, its norm is 1, thus

$$\xi \leq \left\| \mathbf{I} - 2k \phi \right\| \equiv \xi_1$$

Since  $(I - 2k\phi)$  is self-adjoint (see Halmos<sup>(23)</sup> page 180 and Goldstein<sup>(28)</sup> page 24) we may bound  $\xi_1$  as follows

$$\xi_1 = \sup_{\|\underline{x}\|=1} \left| \underline{x}^T (I - 2k\phi) \underline{x} \right| \quad (4.3.1.16)$$

Since  $\phi$  is symmetric positive definite

$$2k\rho_1 \leq 2k\underline{x}^T \phi \underline{x} \leq 2k\rho_N \quad (4.3.1.17)$$

where  $\rho_1$  and  $\rho_N$  are the minimum and maximum eigenvalues of  $\phi$  respectively, and  $\|\underline{x}\| = 1$ . This implies

$$1 - 2k\underline{x}^T \phi \underline{x} \geq 1 - 2k\rho_N \quad (4.3.1.18a)$$

and

$$1 - 2k\underline{x}^T \phi \underline{x} \leq 1 - 2k\rho_1 \quad (4.3.1.18b)$$

thus

$$1 - 2k\rho_N \leq 1 - 2k\underline{x}^T \phi \underline{x} \leq 1 - 2k\rho_1 \quad (4.3.1.19)$$

and

$$\sup_{\|\underline{x}\|=1} \left| \underline{x}^T (I - 2k\phi) \underline{x} \right| = \max \left\{ |1 - 2k\rho_1|, |1 - 2k\rho_N| \right\} \quad (4.3.1.20)$$

Thus

$$\xi \leq \xi_1 = \max \left\{ |1 - 2k\rho_1|, |1 - 2k\rho_N| \right\} \quad (4.3.1.21)$$

If we plot, on  $\xi$  vs  $k$  axes, the two curves  $\xi = |1 - 2k\rho_1|$  and  $\xi = |1 - 2k\rho_N|$  we have

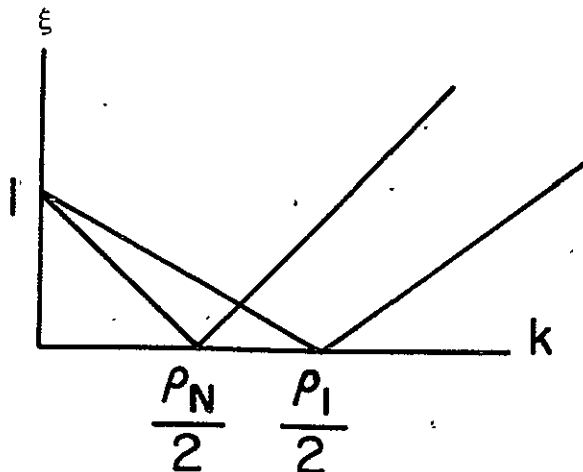


Fig. 4.3.2  $\xi$  vs.  $k$

A plot of  $\xi = \max \left\{ |1 - 2k \rho_1|, |1 - 2k \rho_N| \right\}$  looks like

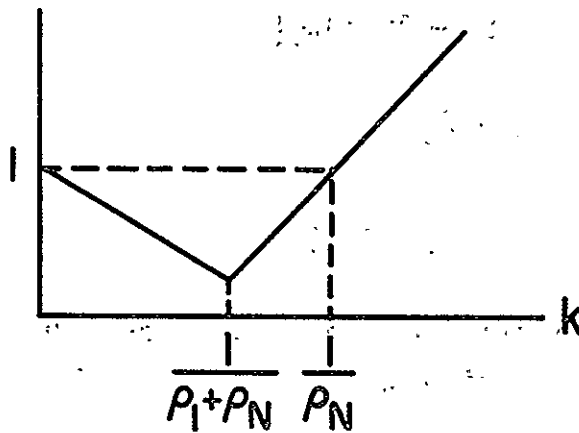


Fig. 4.3.3 Bounds on  $k_{\max}$

The maximum value of  $k$  that still insures convergence ( $k_{\max}$ ) is found by setting

$$-(1 - 2k \rho_N) \leq 1 \quad \text{which implies}$$

$$k \leq \frac{1}{\rho_N} \quad (4.3.1.22)$$

Thus in this section we have proven that our algorithm converges to  $\underline{W}_{\text{opt}}$  for  $k$  sufficiently small. In the next section we will investigate a more useful algorithm, i.e. an algorithm which does not require a priori knowledge of  $\phi$ .

#### Section 4.3.2 The Algorithm, Proof of Convergence, and Bounds on the Rate of Convergence if the Gradient is Estimated

In practice, the mean-square error  $\overline{\epsilon_j^2}$  is normally not available. There are various methods available for estimating  $\overline{\epsilon_j^2}$ . Here we will assume the simplest estimate  $\overline{\epsilon_j^2} \approx \epsilon_j^2$ , i.e. we are approximating the average value of  $\epsilon_j^2$  by its instantaneous value, which is normally available. Thus the  $i^{\text{th}}$  component of the gradient is approximately given by the  $i^{\text{th}}$  partial derivative of  $\epsilon_j^2$  with respect to  $W_i$ .

$$\frac{\partial \epsilon_j^2}{\partial W_i} \approx \frac{\partial \epsilon_j^2}{\partial W_i} = 2 \epsilon_j \frac{\partial \epsilon_j}{\partial W_i}$$

From equation (4.1.1)

$$\frac{\partial \epsilon_j}{\partial W_i} = -s_{ij}$$

thus

$$\nabla (\epsilon_j^2) \approx \nabla \epsilon_j^2 = -2 \epsilon_j \underline{s}_j \quad (4.3.2.1)$$

We will now use this estimated gradient  $\hat{\underline{g}}$  in our algorithm yielding

$$p_{\hat{\underline{g}}} = [I - \underline{n}_1 \underline{n}_1^T] - 2 \epsilon_j \underline{s}_j$$

$$\underline{W}_{j+1} = \underline{W}_j + 2k [I - \underline{n}_1 \underline{n}_1^T] \underline{s}_j \epsilon_j$$

using equation (4.1.1)

$$\underline{W}_{j+1} = \underline{W}_j + 2k [I - \underline{n}_1 \underline{n}_1^T] \underline{s}_j (d_j - \underline{s}_j^T \underline{W}_j) \quad (4.3.2.2)$$

The "initial" condition is

$$\underline{W}_\infty^T \cdot \underline{n}_1 = a \quad (4.3.2.3)$$

$\underline{W}_j$  is now a random vector, and equations (4.3.2.2) represent a set of first-order stochastic difference equations, with forcing stochastic vector  $\underline{s}_j$ .

Let us see what the asymptotic expected value of  $\underline{W}_j$  is:

$$E \{ \underline{W}_{j+1} \} = E \{ \underline{W}_j \} + 2k [I - \underline{n}_1 \underline{n}_1^T] [ \underline{\phi}(\underline{s}, d) - \underline{\phi}(\underline{s}, \underline{s}) E \{ \underline{W}_j \} ]$$



because

$$E \{ \underline{s}_j \underline{s}_j^T \underline{W}_j \} = E \{ \underline{s}_j \underline{s}_j^T \} E \{ \underline{W}_j \}$$

i. e.  $\underline{W}_j$  depends upon  $\underline{s}_1, \dots, \underline{s}_{j-1}$  but is independent of  $\underline{s}_j$ .

From equation (4.1.7a)  $\phi(\underline{s}, d) = \phi(\underline{s}, \underline{s}) \underline{W}_{LMS}$

$$E \{ \underline{W}_{j+1} \} = E \{ \underline{W}_j \} + 2k [ I - \underline{n}_1 \underline{n}_1^T ] \phi [ \underline{W}_{LMS} - E \{ \underline{W}_j \} ] \quad (4.3.2.4)$$

Taking the expected value of (4.3.2.3) yields

$$E \{ \underline{W}_{\infty}^T \} \underline{n}_1 = a \quad (4.3.2.5)$$

Equations (4.3.2.4) and (4.3.2.5) constitute a set of first order deterministic difference equations, exactly the same as equations (4.3.1.4)

and (4.3.1.5). Thus the solution (unique since  $\text{rank} (I - \underline{n}_1 \underline{n}_1^T) \phi = n-1$ ) is

$$E \{ \underline{W}_{\infty} \} = \underline{W}_{\text{optimum}} = \frac{(a - \underline{W}_{LMS}^T \underline{n}_1)}{(\underline{n}_1^T \phi^{-1} \underline{n}_1)} \phi^{-1} \underline{n}_1 + \underline{W}_{LMS} \quad (4.3.2.6)$$

We have shown that the mean of  $\underline{W}_j$  converges to  $\underline{W}_{\text{opt}}$ . However, since equations (4.3.2.2) are stochastic, we must also investigate the behavior of the variance of the random weight vector  $\underline{W}_j$  about its expected asymptotic value, given by  $E \{ \underline{W}_{\infty} \} = \underline{W}_{\text{optimum}} \equiv \underline{\overline{W}}_{\infty}$ .

$$\text{Let } \underline{q}_j \equiv \underline{W}_j - \underline{\overline{W}}_{\infty} \quad (4.3.2.7)$$

In terms of  $\underline{q}$ , the algorithm (4.3.2.2) becomes

$$\begin{aligned}\underline{q}_{j+1} = & \underline{q}_j - k \cdot 2 [I - \underline{n}_1 \underline{n}_1^T] \underline{s}_j \underline{s}_j^T \underline{q}_j \\ & - k \cdot 2 [I - \underline{n}_1 \underline{n}_1^T] \underline{s}_j \underline{s}_j^T \underline{\bar{W}}_\infty \\ & + k \cdot 2 [I - \underline{n}_1 \underline{n}_1^T] \underline{s}_j d_j\end{aligned}$$

$$\text{Define } T_j \equiv 2 \underline{s}_j \underline{s}_j^T \quad (4.3.2.8)$$

$$\underline{V}_j \equiv 2 d_j \underline{s}_j \quad (4.3.2.9)$$

$$H_j \equiv (I - \underline{n}_1 \underline{n}_1^T) T_j \quad (4.3.2.10)$$

thus

$$\underline{q}_{j+1} = \underline{q}_j - k [I - \underline{n}_1 \underline{n}_1^T] [T_j \underline{q}_j + T_j \underline{\bar{W}}_\infty - \underline{V}_j] \quad (4.3.2.11)$$

This may be rewritten as

$$\underline{q}_{j+1} = \underline{q}_j - k \underline{\varphi}_j \quad (4.3.2.12)$$

where

$$\underline{\varphi}_j \equiv H_j \underline{q}_j + \underline{h}_j \quad (4.3.2.13)$$

$$\underline{h}_j \equiv (I - \underline{n}_1 \underline{n}_1^T) (T_j \underline{\bar{W}}_\infty - \underline{V}_j) \quad (4.3.2.14)$$

Note that  $E\{H_j\}$  and  $E\{\underline{h}_j\}$  are independent of  $j$ . Also  $H_j$  and  $\underline{h}_j$  are statistically independent of  $H_k$  and  $\underline{h}_k$  if  $j \neq k$ , because we assumed that  $\underline{s}_j$  and  $\underline{s}_k$  are statistically independent for  $k \neq j$ .

Noting that

$$E\{T_j\} = 2 \phi(\underline{s}, \underline{s})$$

and

$$E \{ \underline{V}_j \} = 2 \underline{\phi}(\underline{s}, d) = 2 \underline{\phi}(\underline{s}, \underline{s}) \underline{W}_{LMS}$$

it is easily shown that

$$E \{ \underline{h}_j \} = \underline{0} \quad (4.3.2.15)$$

$$\text{Note that } E \{ H_j \} = 2 (I - \underline{n}_1 \underline{n}_1^T) \underline{\phi}(\underline{s}, \underline{s}) \equiv \underline{0} \quad (4.3.2.16)$$

The algorithm is thus

$$\underline{q}_{j+1} = \underline{q}_j - k \underline{\varphi}_j$$

where

$$\underline{\varphi}_j = H_j \underline{q}_j + \underline{h}_j$$

$H_j$  is a sequence of random  $n \times n$  matrices;  $\underline{h}_j$  is a sequence of random  $n$ -tuple vectors; the expected values of  $H_j$  and  $\underline{h}_j$  were shown to be independent of  $j$ ;  $H_j$  and  $\underline{h}_j$  are independent of  $H_l$  and  $\underline{h}_l$  for  $l \neq j$ ;  $E\{\underline{h}_j\} = \underline{0}$ ; and the elements of  $H_j$  and  $\underline{h}_j$  have finite variance, with  $E\{H_j\} = \underline{0}$

Under these conditions, it is shown in appendix A that for  $k$  sufficiently small,

$$\lim_{j \rightarrow \infty} || E \{ \underline{q}_j \} || = 0 \quad (4.3.2.17)$$

$$\text{and } \lim_{j \rightarrow \infty} \sup || \underline{q}_j || \leq V(k) \quad (4.3.2.18)$$

where the norm of a random vector  $\underline{u}$  is defined as

$$|| \underline{u} || \equiv \sqrt{E \{ \underline{u}^T \underline{u} \}} \quad (4.3.2.19)$$

and

$$\lim_{k \rightarrow 0} V(k) = 0 \quad (4.3.2.20)$$

Equation (4.3.2.17) shows again that the weights converge to  $\underline{W}_{\text{optimum}}$  and equation (4.3.2.18) shows that the variance of the random weight vector about its expected value is bounded, and the bound can be made as small as desired by choosing  $k$  sufficiently small.

The rate of convergence of the mean of the weight vector is shown in the proof of the above theorem to be bounded by  $\xi$ , where

$$\xi = || I - k (I - \underline{n}_1 \underline{n}_1^T) ||_2 \phi || \quad (4.3.2.21)$$

and  $0 < \xi < 1$  as shown in section 4.3.1.

### Section 4.3.3      The Algorithm, proof of Convergence, and Bounds on the Rate of Convergence if the Gradient is Estimated, and the Estimate is Noisy.

When our estimated gradient contains noise, wherever we have the quantity  $\underline{s}_j$  in section 4.3.2 we replace it by  $\underline{s}_j + \underline{n}_j$ . To characterize the noise we will assume

$E\{\underline{n}_j\} = \underline{0}$ ,  $E\{\underline{n}_j \underline{n}_j^T\} = \phi_n$ , and  $\underline{s}_j, \underline{s}_k, \underline{n}_l, \underline{n}_m$  are statistically independent for  $k \neq j$  and  $n \neq m$ .

The algorithm becomes

$$\underline{W}_{j+1} = \underline{W}_j + 2k [I - \underline{n}_1 \underline{n}_1^T] (\underline{s}_j + \underline{n}_j) [d_j - (\underline{s}_j^T + \underline{n}_j^T) \underline{W}_j] \quad (4.3.3.1)$$

with the initial condition the same as before, i.e.

$$\underline{W}_{\infty}^T \underline{n}_1 = a \quad (4.3.3.2)$$

Equations (4.3.3.1) represent a set of first-order stochastic difference equations, with forcing stochastic vectors  $\underline{s}_j$  and  $\underline{n}_j$ .

Let's find the asymptotic expected value of  $\underline{W}_j$

$$E\{\underline{W}_{j+1}\} = E\{\underline{W}_j\} + 2k [I - \underline{n}_1 \underline{n}_1^T] [\phi(\underline{s}, d) - \phi(\underline{s}, \underline{s}) E\{\underline{W}_j\} - \phi_n E\{\underline{W}_j\}]$$

Using (4.1.7a) and setting  $E\{\underline{W}_{j+1}\} = E\{\underline{W}_j\} = \underline{\bar{W}}_\infty$  gives

$$\underline{0} = \left[ \underline{I} - \underline{n}_1 \underline{n}_1^T \right] \phi \left[ (\underline{I} + \phi^{-1} \phi_n) \underline{\bar{W}}_\infty - \underline{W}_{LMS} \right] \quad (4.3.3.3)$$

Taking the expected value of (4.3.3.2) yields

$$\underline{\bar{W}}_\infty^T \underline{n}_1 = a \quad (4.3.3.4)$$

$$\text{Define } \underline{x} \triangleq \left[ (\underline{I} + \phi^{-1} \phi_n) \underline{\bar{W}}_\infty - \underline{W}_{LMS} \right] \quad (4.3.3.5)$$

$$\text{then } \left[ \underline{I} - \underline{n}_1 \underline{n}_1^T \right] \phi \underline{x} = \underline{0} \quad (4.3.3.6)$$

By the previous arguments, a solution to (4.3.3.6) exists and is unique because

$$\text{rank} \left[ \underline{I} - \underline{n}_1 \underline{n}_1^T \right] \phi = n-1$$

Equation (4.3.3.6) is the same as equation (4.3.1.7), thus the solution is given by

$$\underline{x} = a \phi^{-1} \underline{n}_1 \quad (4.3.3.7)$$

where the value of  $a$  is chosen so as to satisfy the initial condition, given by equation (4.3.3.4), i.e.

$$a = \frac{\left[ a - \underline{n}_1^T (\underline{I} + \phi^{-1} \phi_n)^{-1} \underline{W}_{LMS} \right]}{\underline{n}_1^T (\underline{I} + \phi^{-1} \phi_n)^{-1} \phi^{-1} \underline{n}_1} \quad (4.3.3.8)$$

The solution for  $\underline{\bar{W}}_\infty$  is thus

$$\underline{\bar{W}}_\infty = (\underline{I} + \phi^{-1} \phi_n)^{-1} \left\{ \underline{W}_{LMS} + \frac{\left[ a - \underline{n}_1^T (\underline{I} + \phi^{-1} \phi_n)^{-1} \underline{W}_{LMS} \right] \phi^{-1} \underline{n}_1}{\underline{n}_1^T (\underline{I} + \phi^{-1} \phi_n)^{-1} \phi^{-1} \underline{n}_1} \right\} \quad (4.3.3.9)$$

Remembering that

$$\underline{W}_{\text{opt}} = \underline{W}_{\text{LMS}} + \frac{(\underline{a} - \underline{n}_1^T \underline{W}_{\text{LMS}})}{(\underline{n}_1^T \phi^{-1} \underline{n}_1)} \phi^{-1} \underline{n}_1$$

we see that  $\bar{\underline{W}}_{\infty}$  differs from  $\underline{W}_{\text{opt}}$  in this case, i.e. a bias exists, and the bias approaches zero as the noise matrix  $\phi_n$  approaches the zero matrix.

Again, since the weight vectors are random, before we can conclude that the weight vectors converge to  $\underline{W}_{\text{optimum}}$ , we must examine the variations of the weight vectors about their asymptotic expected value, given by (4.3.3.8)

$$\text{Define } \underline{q}_j \equiv \underline{W}_j - \bar{\underline{W}}_{\infty} \quad (4.3.3.10)$$

In terms of  $\underline{q}$ , the algorithm (4.3.3.1) becomes

$$\begin{aligned} \underline{q}_{j+1} = & \underline{q}_j - k_2 [I - \underline{n}_1 \underline{n}_1^T] (\underline{s}_j + \underline{n}_j) (\underline{s}_j^T + \underline{n}_j^T) \underline{q}_j \\ & - k_2 [I - \underline{n}_1 \underline{n}_1^T] (\underline{s}_j + \underline{n}_j) (\underline{s}_j^T + \underline{n}_j^T) \bar{\underline{W}}_{\infty} \\ & + k_2 [I - \underline{n}_1 \underline{n}_1^T] d_j (\underline{s}_j + \underline{n}_j) \end{aligned}$$

$$\text{Define } \underline{T}_j \equiv 2(\underline{s}_j + \underline{n}_j) (\underline{s}_j^T + \underline{n}_j^T) \quad (4.3.3.11)$$

$$\underline{V}_j \equiv 2 d_j (\underline{s}_j + \underline{n}_j) \quad (4.3.3.12)$$

$$\underline{H}_j \equiv (I - \underline{n}_1 \underline{n}_1^T) \underline{T}_j \quad (4.3.3.13)$$

$$\underline{q}_{j+1} = \underline{q}_j - k [I - \underline{n}_1 \underline{n}_1^T] [\underline{T}_j \underline{q}_j + \underline{T}_j \bar{\underline{W}}_{\infty} - \underline{V}_j] \quad (4.3.3.14)$$

This may be rewritten as

$$\underline{q}_{j+1} = \underline{q}_j - k \underline{\phi}_j \quad (4.3.3.15)$$

$$\text{where } \underline{\phi}_j \equiv H_j \underline{q}_j + \underline{h}_j \quad (4.3.3.16)$$

$$\underline{h}_j \equiv (I - \underline{n}_1 \underline{n}_1^T) (T_j \underline{\overline{W}}_\infty - \underline{V}_j) \quad (4.3.3.17)$$

Note that  $E \{ H_j \}$  and  $E \{ \underline{h}_j \}$  are independent of  $j$ . Also  $H_j$  and  $\underline{h}_j$  are statistically independent of  $H_k$  and  $\underline{h}_k$  if  $j \neq k$  because we assumed  $\underline{s}_j, \underline{s}_k, \underline{n}_l, \underline{n}_m$  are statistically independent for  $k \neq j$  and  $n \neq m$ .

Noting that

$$E \{ H_j \} = 2 [ \phi + \phi_n ]$$

$$E \{ \underline{V}_j \} = 2 \underline{\phi}(\underline{s}, d) = 2 \phi \underline{W}_{LMS}$$

it is easy to show that

$$E \{ \underline{h}_j \} = \underline{0} \quad (4.3.3.18)$$

$$\text{Also } E \{ H_j \} = 2 (I - \underline{n}_1 \underline{n}_1^T) (\phi + \phi_n) \equiv \underline{a} \quad (4.3.3.19)$$

By the same argument as before, we may show that for  $k$  sufficiently small,

$$\lim_{j \rightarrow \infty} || E \{ \underline{q}_j \} || = 0 \quad (4.3.3.20)$$

and

$$\limsup_{j \rightarrow \infty} || \underline{q}_j || \leq V(k) \quad (4.3.3.21)$$

This proves convergence.

Again the rate of convergence is bounded by  $\xi$ , which depends upon  $k$ , the eigenvalues of  $(\phi + \phi_n)$ , and the constraint.

#### Section 4.4 Simulation Results

As a check on the theoretical work we have done in this chapter, we programmed the following algorithms on IBM 360/50 in Fortran IV.

Let us first consider the algorithm given by equation (4.3.1.4) where the gradient is assumed to be known. We let the dimension of the vector  $\underline{W}$  be four.

Let

$$d_j \equiv [1 \ 1 \ 1 \ 1] \begin{bmatrix} s_{1j} \\ s_{2j} \\ s_{3j} \\ s_{4j} \end{bmatrix} \quad (4.4.1)$$

$$\phi(\underline{s}, \underline{s}) \equiv \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 4 \end{bmatrix} \quad (4.4.2)$$

i. e. all components of the vector  $\underline{s}_j$  are assumed to be gaussian, zero mean, and uncorrelated.

Thus

$$\phi^{-1}(\underline{s}, \underline{s}) = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & \frac{1}{2} & 0 & 0 \\ 0 & 0 & \frac{1}{3} & 0 \\ 0 & 0 & 0 & \frac{1}{4} \end{bmatrix} \quad (4.4.3)$$

and

$$\phi(\underline{s}, \underline{d}) = \begin{bmatrix} E\{s_{1j}d_j\} \\ E\{s_{2j}d_j\} \\ E\{s_{3j}d_j\} \\ E\{s_{4j}d_j\} \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} \quad (4.4.4)$$



The LMS value of the weights is given by

$$\underline{W}_{LMS} = \phi^{-1}(\underline{s}, \underline{s}) \phi(\underline{s}, d) = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad (4.4.5)$$

For our constraint we let

$$\underline{n}_1 = \begin{bmatrix} \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} \\ 0 \\ 0 \end{bmatrix} \quad \text{and } a = 3 \quad (4.4.6)$$

i.e. the linear constraint equation is

$$W_1 + W_2 \geq 3\sqrt{2}$$

which means that there are no constraints on  $W_3$  and  $W_4$ .

For our initial conditions we considered two cases:

$$\underline{W}_0 = \begin{bmatrix} 3\sqrt{2} \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \text{or } \underline{W}_0 = \begin{bmatrix} 10+3\sqrt{2} \\ 10 \\ 0 \\ 0 \end{bmatrix}$$

which exactly satisfy the constraint.

The Lagrange solution is

$$\underline{W}_{opt} = \begin{bmatrix} 1 + 2\sqrt{2} \\ 1 - 2\sqrt{2} \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3.82 \\ -0.41 \\ 1.0 \\ 1.0 \end{bmatrix} \quad (4.4.7)$$

A limit on the values of  $k$  which insure convergence is, from equation (4.3.1.22), given by those values of  $k$  for which  $\xi < 1$ , i.e.

$$k \leq \frac{1}{4}$$

The algorithm is (see equation 4.3.1.4).

$$\underline{W}_{j+1} = \underline{W}_j - 2k [I - \underline{n}_1 \underline{n}_1^T] \phi [\underline{W}_j - \underline{W}_{LMS}] \quad (4.4.8)$$

where

$$\underline{W}_0 = \begin{bmatrix} 3\sqrt{2} \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \underline{W}_0 = \begin{bmatrix} 10 + 3\sqrt{2} \\ 10 \\ 0 \\ 0 \end{bmatrix}$$

Using the values we have chosen for  $\underline{n}_1$  and  $\phi$ , the algorithm may be rewritten as

$$\underline{W}_{j+1} = \underline{W}_j - k \begin{bmatrix} W_{1j} + 2W_{2j} - 3 \\ W_{1j} + 2W_{2j} - 3 \\ 6W_{3j} - 6 \\ 8W_{4j} - 8 \end{bmatrix} \quad (4.4.9)$$

In the steady state,  $\underline{W}$  should converge to  $\underline{W}_\infty =$   
and the asymptotic MSE should be given by

$$\begin{bmatrix} 3.82 \\ -0.41 \\ 1.0 \\ 1.0 \end{bmatrix}$$

$$E \{ \epsilon_j^2 \} = \overline{d_j^2} - 2 \underline{\phi}^T (\underline{s}, d) \underline{W} + \underline{W}^T \underline{\phi} \underline{W}$$

evaluated at  $\underline{W} = \underline{W}_\infty \approx \underline{W}_{opt}$  which is

$$E \{ \epsilon_j^2 \} \Big|_{\underline{W} = \underline{W}_\infty = \underline{W}_{opt}} = 12.0 \quad (4.4.10)$$

We ran the above algorithm for various values of  $k$ , with the initial condition  $\underline{W}_0 = \text{col} [10 + 3\sqrt{2}, 10, 0, 0]$  and the results are shown in Fig 4.4.1. Note that as  $k$  increased from 0.01 to 0.25 (above which we no longer have convergence, theoretically or in the simulation, as demonstrated by Fig 4.4.1 when we let  $k=0.252$ ) the rate of convergence agrees with the bound given by Fig. 4.3.3. Fig. 4.4.2. shows how the norm of the vector  $\underline{q}$  (see equation 4.3.1.9) converges to zero for various values of  $k$ . From this graph we can compare the actual time constant, for a particular value of  $k$ , to the theoretical bound on the time constant ( $\xi$ ), e.g. for  $k = 0.01$ ,  $|\underline{q}|$  decreased from 15.82 to 12.26 in ten iterations. Setting  $12.26 = \xi^{10} (15.82)$  implies  $\xi = .975$  which is in agreement with Fig. 4.3.3, which bounds the rate of convergence for this value of  $k$  by  $1 - 2k\rho_1 = 0.98$ .

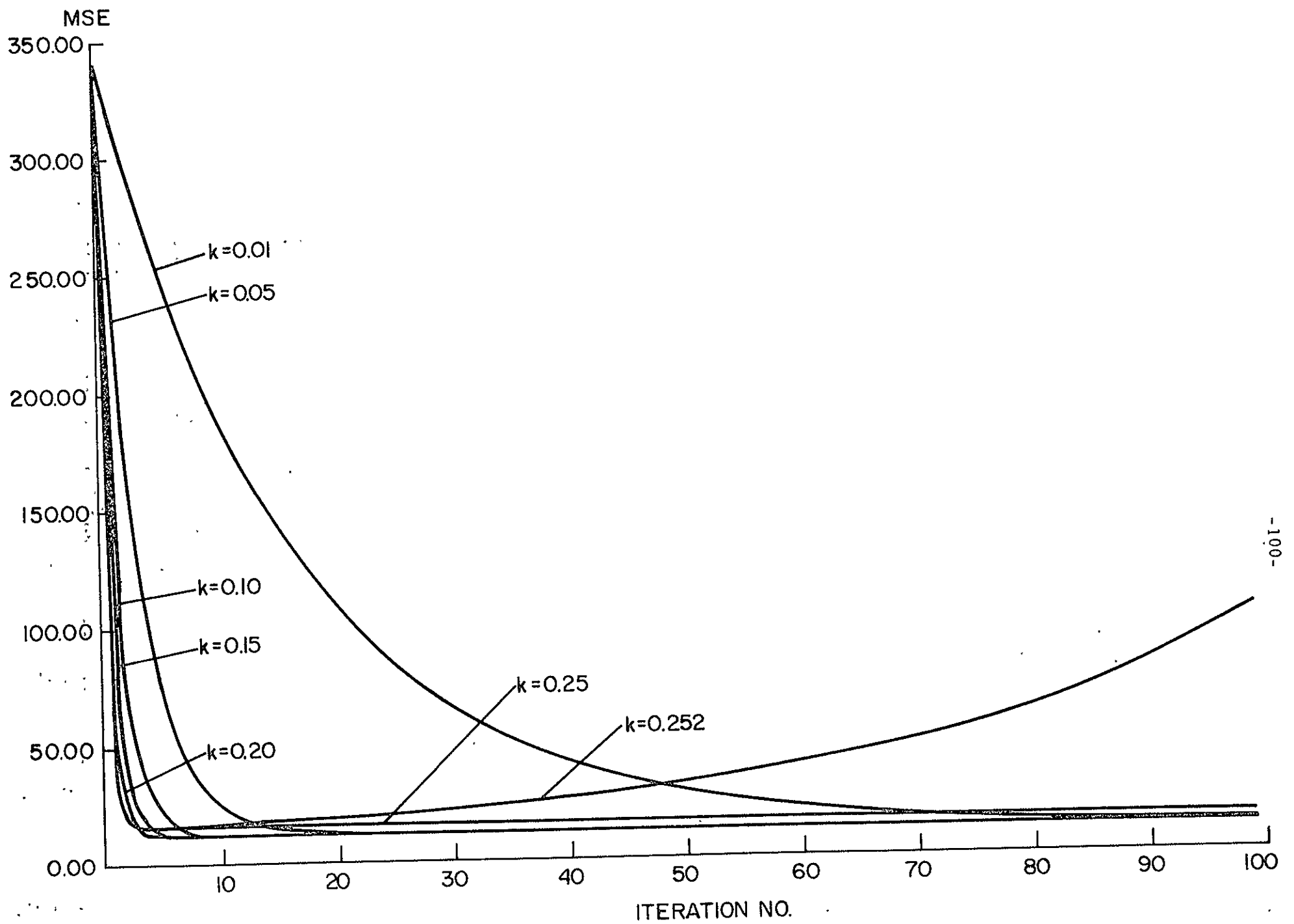


Fig. 4.4.1 Gradient Known , No Additive Noise

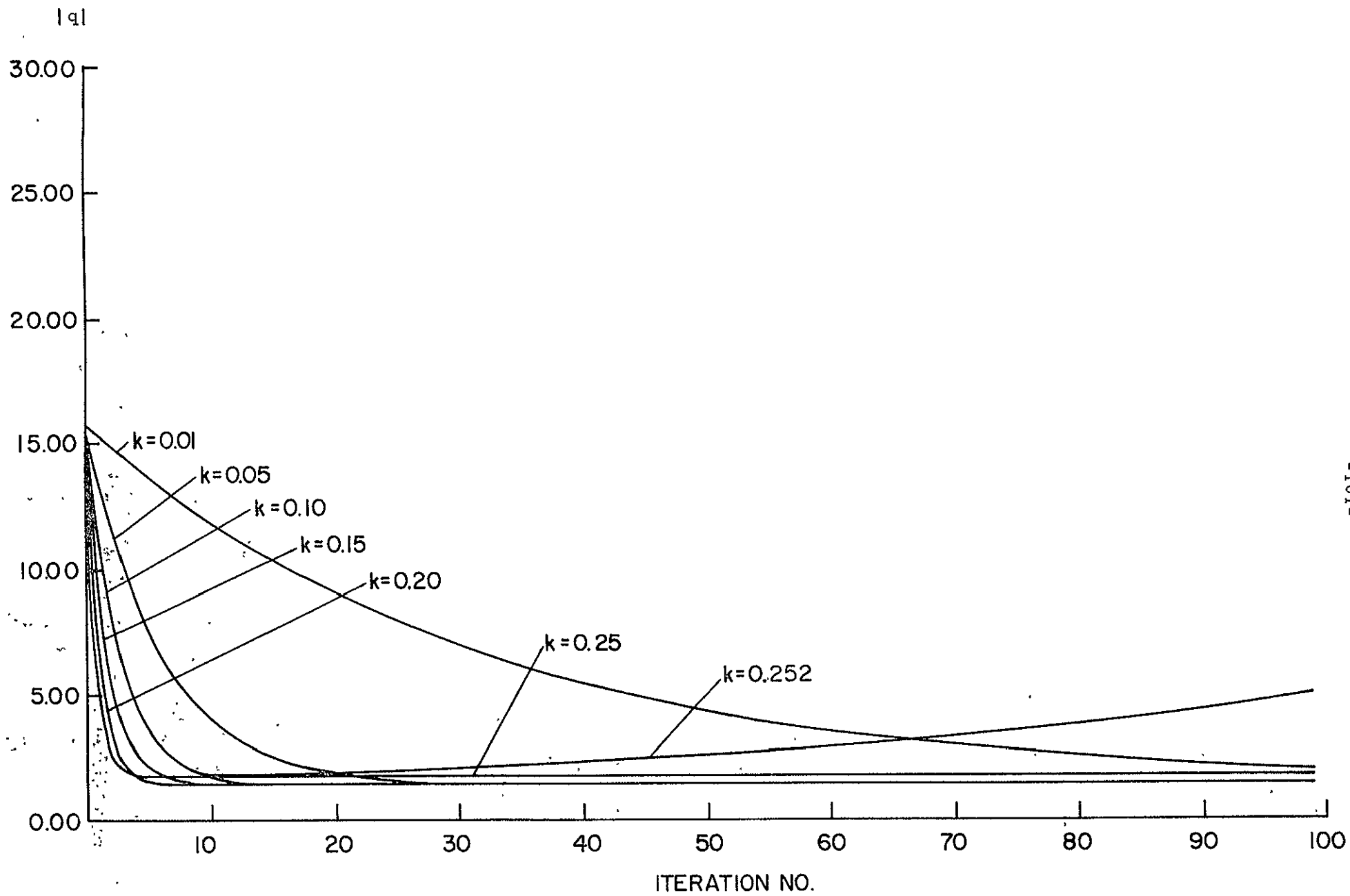


Fig. 4.4.2 Gradient Known , No Additive Noise

If the gradient must be estimated from the incoming data, the algorithm (see equation 4.3.2.2) is

$$\underline{W}_{j+1} = \underline{W}_j + 2k [I - \underline{n}_1 \underline{n}_1^T] \underline{s}_j [d_j - \underline{s}_j^T \cdot \underline{W}_j] \quad (4.4.11)$$

Using the values we have chosen for  $\underline{n}_1$  and  $d_j$  the algorithm may be rewritten in the form

$$\underline{W}_{j+1} = \underline{W}_j + k \begin{bmatrix} (s_{j1} + s_{j2}) u_j \\ (s_{j1} + s_{j2}) u_j \\ 2 s_{j3} u_j \\ 2 s_{j4} u_j \end{bmatrix} \quad (4.4.12)$$

$$\text{where } u_j = (s_{j1} + s_{j2} + s_{j3} + s_{j4}) - (s_{j1} W_{j1} + s_{j2} W_{j2} + s_{j3} W_{j3} + s_{j4} W_{j4}) \quad (4.4.13)$$

In the steady-state,  $\underline{W}_j$  should converge to the same values as before, and the asymptotic MSE should be 12.0.

The results of the simulation for  $k = 0.01$  and  $\underline{W}_0 =$

$$\begin{bmatrix} 10 + 3\sqrt{2} \\ 10 \\ 0 \\ 0 \end{bmatrix}$$

are shown in Figs. 4.4.3 and 4.4.4 and agree with the theoretical values above.

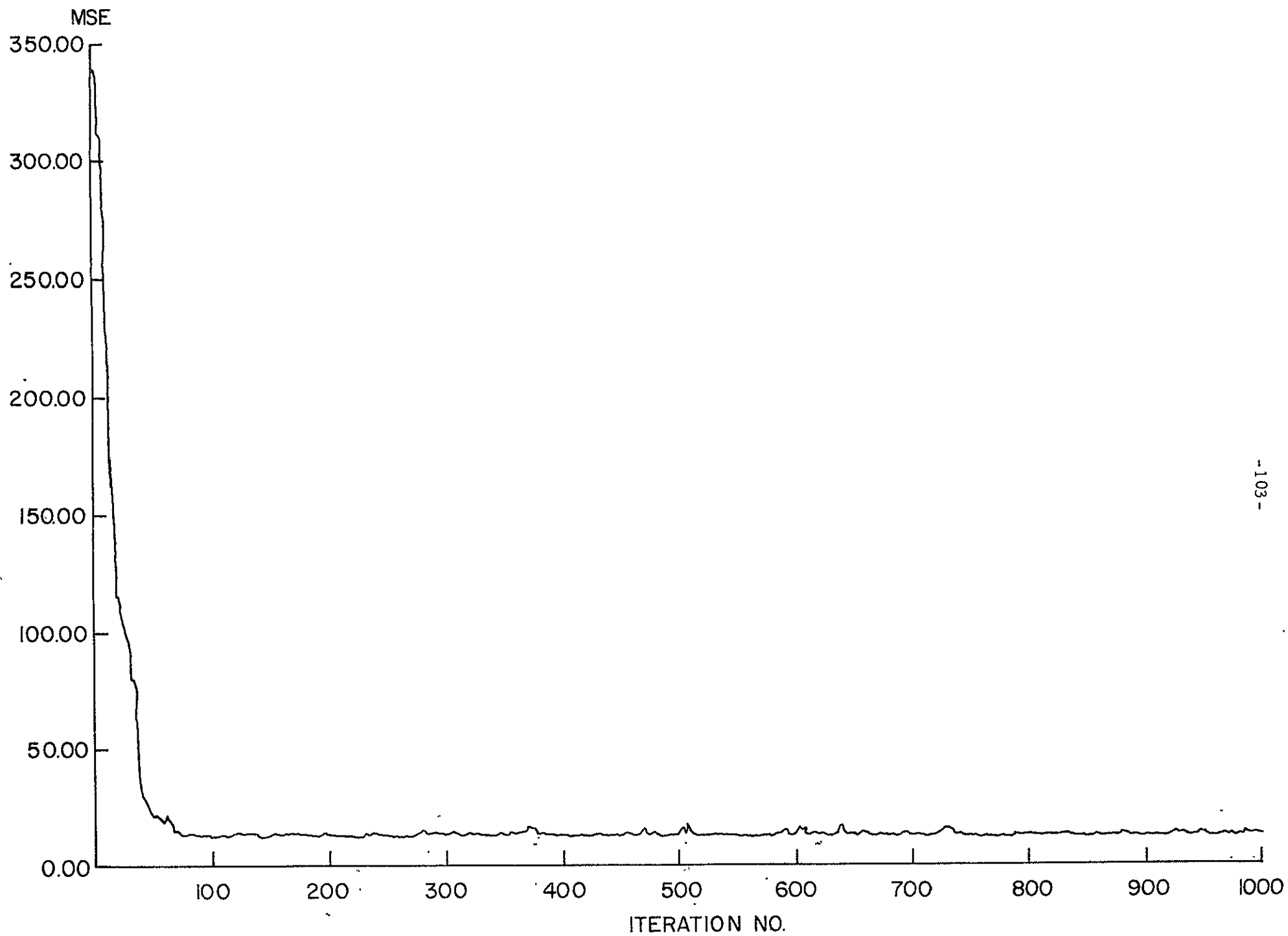


Fig. 4.4.3 Gradient Estimated , No Additive Noise

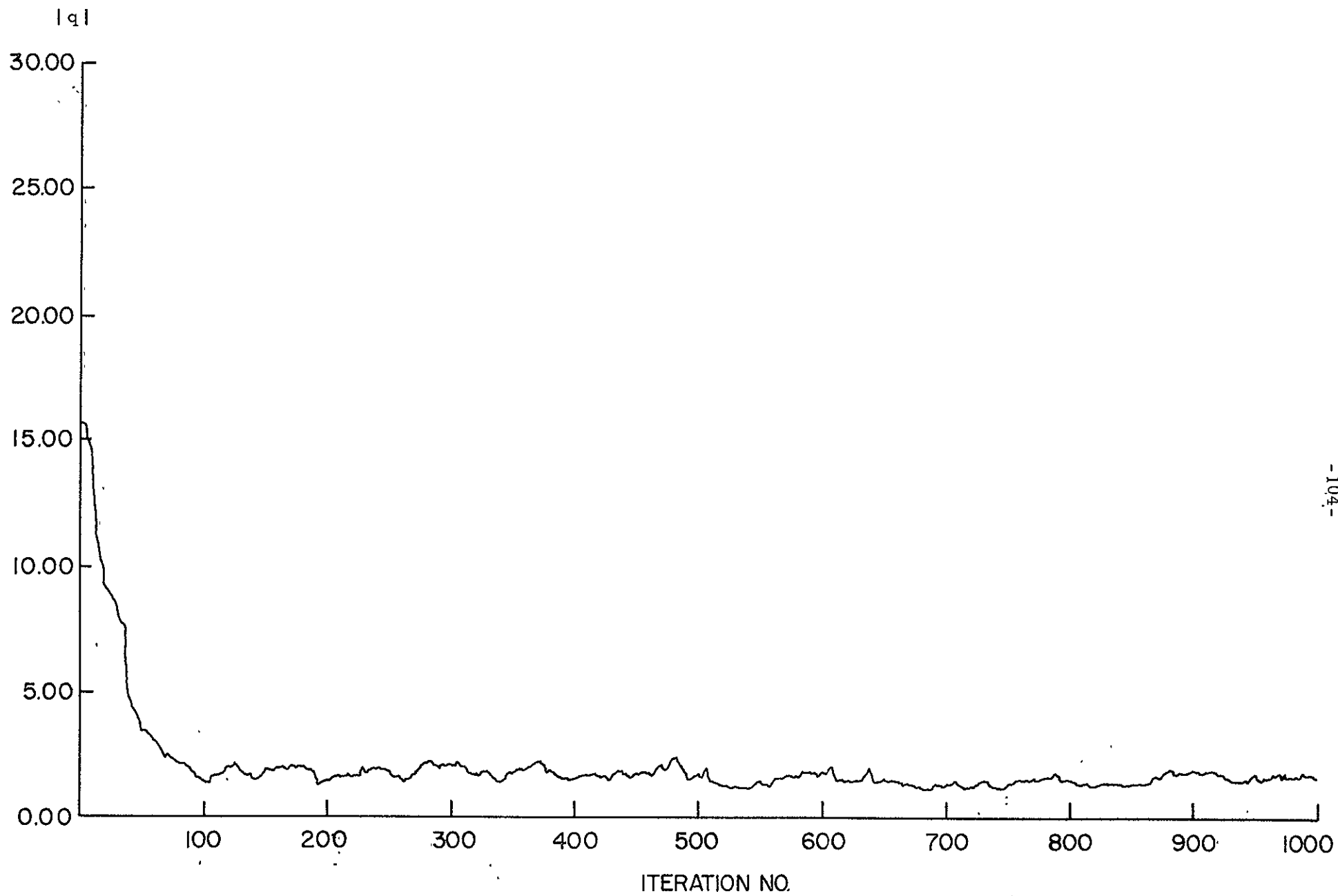


Fig. 4.4.4 Gradient Estimated , No Additive Noise



Finally if the gradient must be estimated from the incoming data, and the incoming data is noisy, the algorithm (see equation 4.3.3.1) becomes

$$\underline{W}_{j+1} = \underline{W}_j + 2k [I - \underline{n}_1 \underline{n}_1^T] (\underline{s}_j + \underline{n}_j) [d_j - (\underline{s}_j^T + \underline{n}_j^T) \underline{W}_j] \quad (4.4.14)$$

Using our specific values for the above quantities, the algorithm may be rewritten as

$$\underline{W}_{j+1} = \underline{W}_j + k u_j \begin{bmatrix} s_{j1} + n_{j1} + s_{j2} + n_{j2} \\ s_{j1} + n_{j1} + s_{j2} + n_{j2} \\ 2(s_{j3} + n_{j3}) \\ 2(s_{j4} + n_{j4}) \end{bmatrix} \quad (4.4.15)$$

where

$$u_j \equiv d_j - (\underline{s}_j^T + \underline{n}_j^T) \underline{W}_j \quad (4.4.16)$$

Let the noise correlation matrix be

$$\phi_n = 0.1 I \quad (4.4.17)$$

In this case, see equation (4.3.3.9), the average asymptotic value  $\underline{\bar{W}}_\infty$  should be

$$\underline{\bar{W}}_\infty = \begin{bmatrix} 3.72 \\ -0.515 \\ 0.967 \\ 0.975 \end{bmatrix} \quad (4.4.18)$$

and the asymptotic MSE should be  $\approx 11.9$ . The results of this

simulation, for  $k=0.01$  and  $\underline{W}_0 = \begin{bmatrix} 10 + 3\sqrt{2} \\ 10 \\ 0 \\ 0 \end{bmatrix}$  are shown in Figs. 4.4.5 and 4.4.6

Figs. 4.4.7 - 4.4.10 indicate how the convergence rate and asymptotic MSE change as the additive noise in the incoming data increases. Figs. 4.4.7 and 4.4.8 correspond to  $\phi_n = 1.0 I$ ,  $\underline{W}_\infty = \text{col} [ 3.17 \quad -1.12 \quad 0.747 \quad 0.80 ]$ , and asymptotic MSE 14.0. Figs. 4.4.9 and 4.4.10 correspond to  $\phi_n = 10.0 I$ ,  $\underline{W}_\infty = \text{col} [ 2.34 \quad -1.90 \quad 0.28 \quad 0.29 ]$ , and asymptotic MSE  $\approx 28$ .

Comparing Figs. 4.4.5, 4.4.7, 4.4.9, and 4.4.3 we see that it took longer to converge when we had additive noise than when we did not have additive noise in the incoming data.

In Figs. 4.4.11 and 4.4.12 we kept everything the same as in Figs. 4.4.5 and 4.4.6 except that we started at  $\underline{W}_0 = \text{col} [ 3\sqrt{2}, 0, 0, 0 ]$  which is much closer to the steady-state value,  $\underline{W}_\infty$ , and expanded the vertical scale. From these figures we notice that the MSE is somewhat sensitive to the occasional noise sample whose value is greater than three or four standard deviations away from the mean value of the noise which in our case is zero. This suggests that one might achieve a smaller value for the steady state variance if the algorithm were

$$\underline{W}_{j+1} = \underline{W}_j - k f [ \nabla (\text{MSE}) ]$$

$$\text{where } f [ \nabla (\text{MSE}) ] = \begin{cases} \nabla (\text{MSE}) & \text{if } \nabla (\text{MSE}) \leq K_0 \\ K_0 & \text{if } \nabla (\text{MSE}) > K_0 \end{cases}$$

However, this approach was not investigated further.

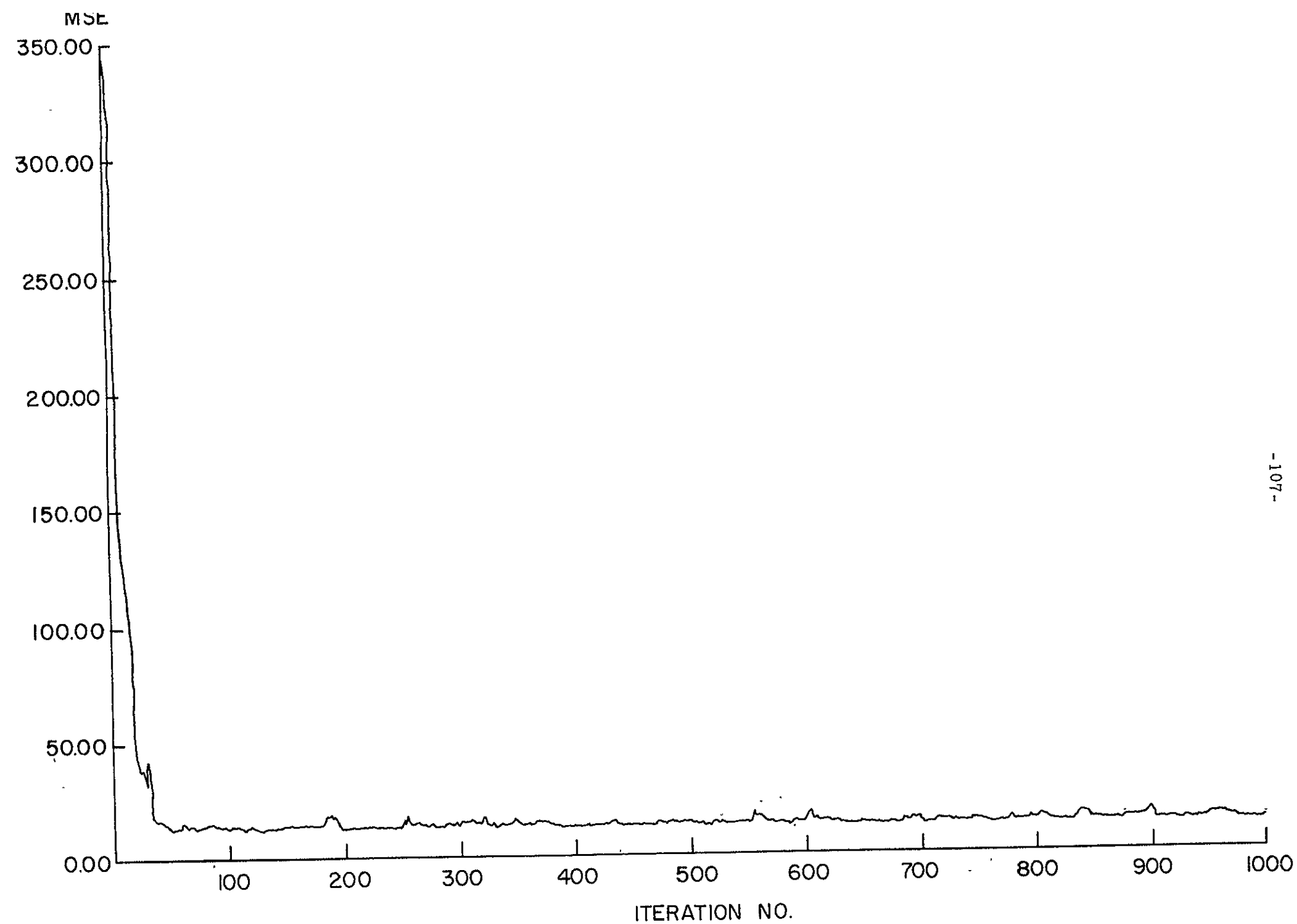


Fig. 4.4.5 Gradient Estimated, Plus Additive Noise

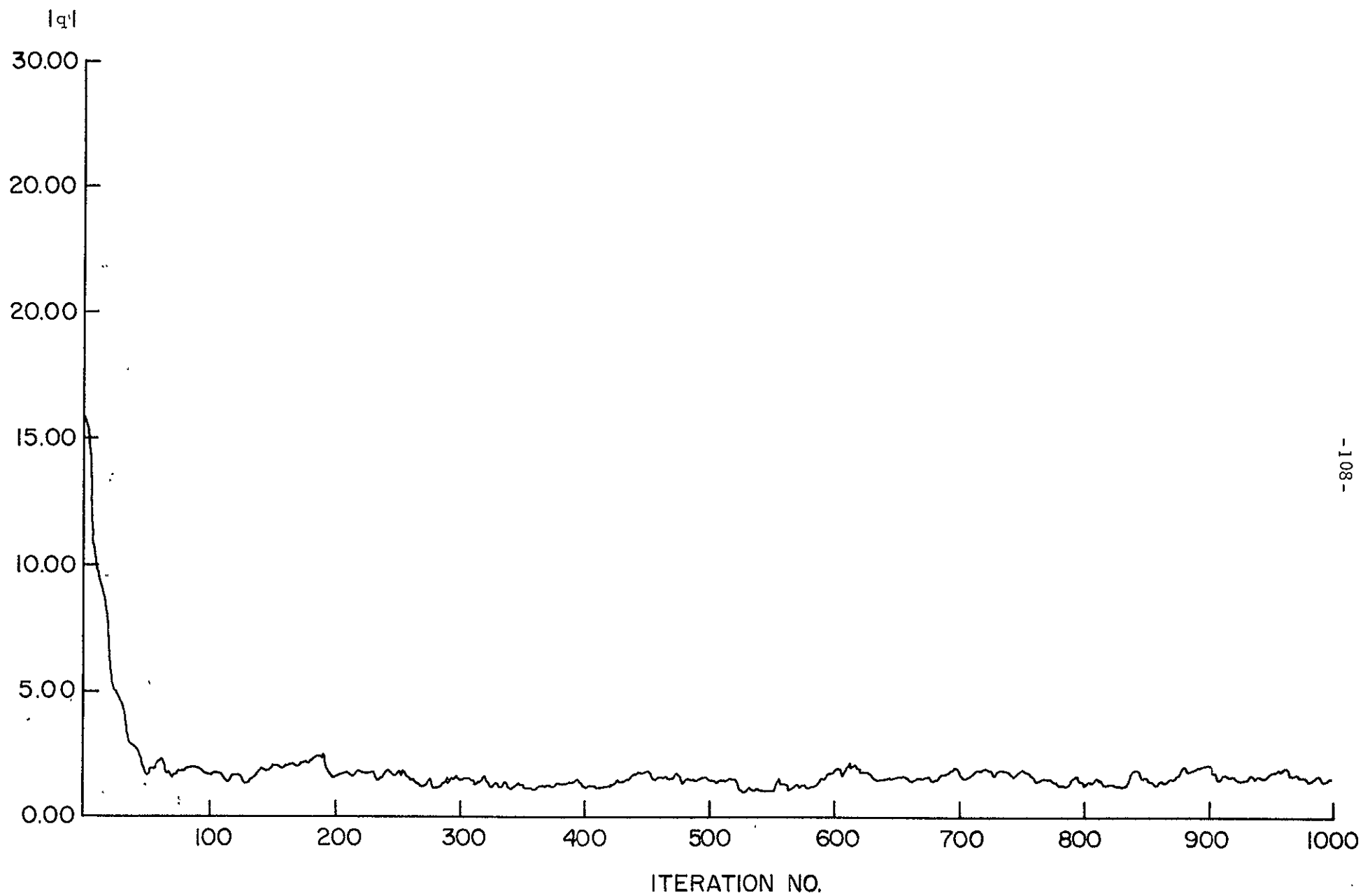


Fig. 4.4.6 Gradient Estimated , Plus Additive Noise

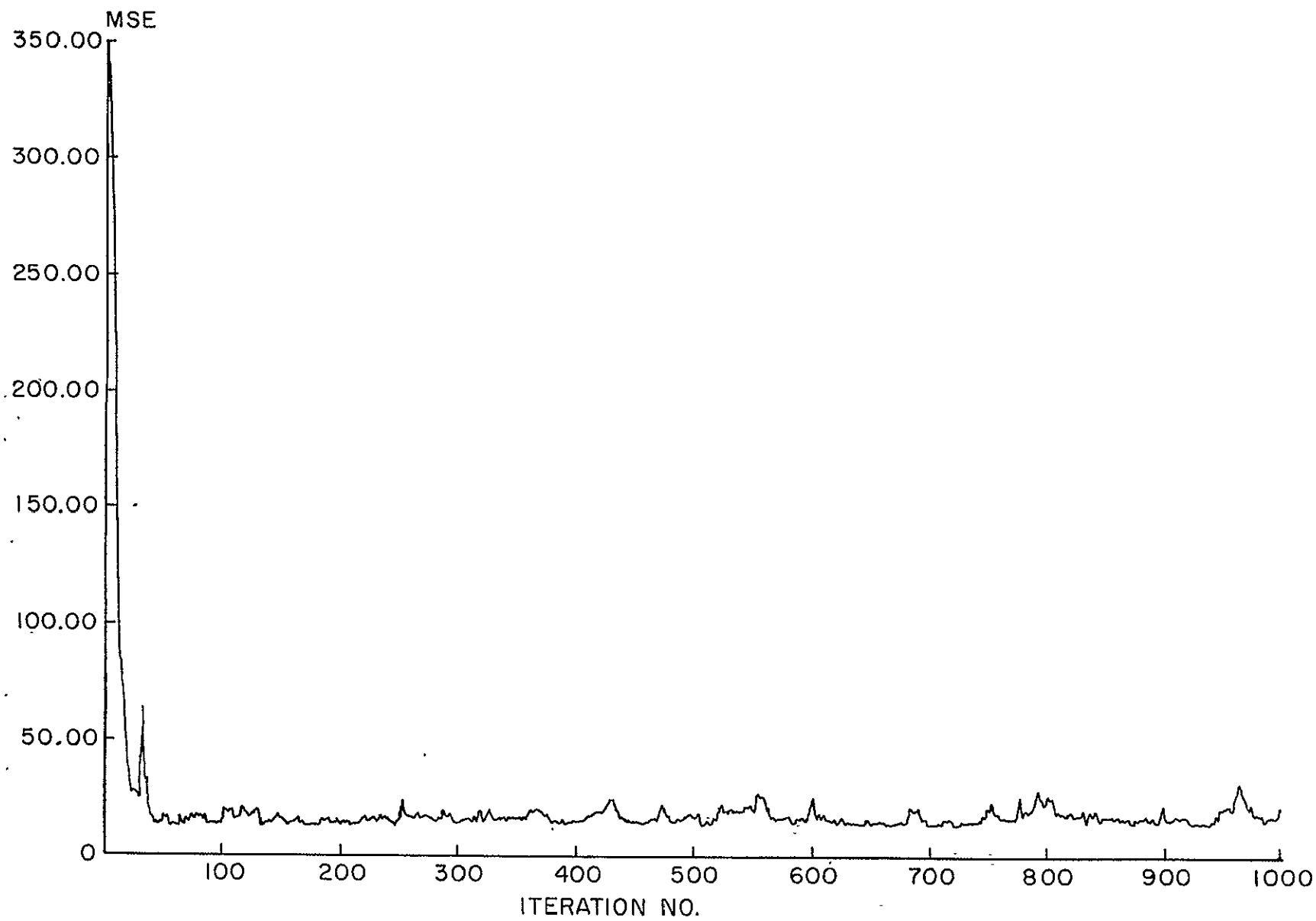


Fig. 4.4.7 Gradient Estimated, Plus Additive Noise

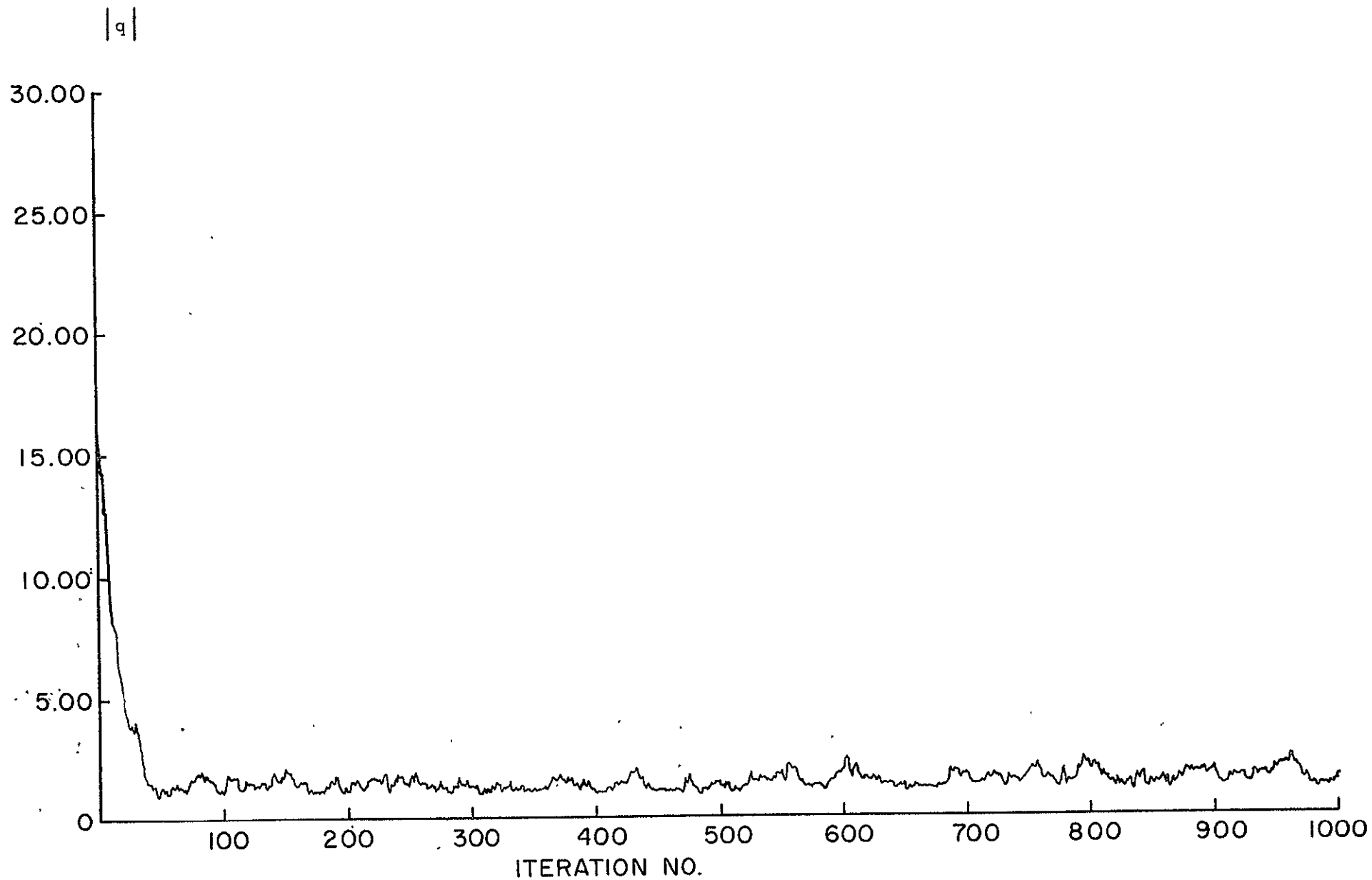


Fig. 4.4.8 Gradient Estimated, Plus Additive Noise

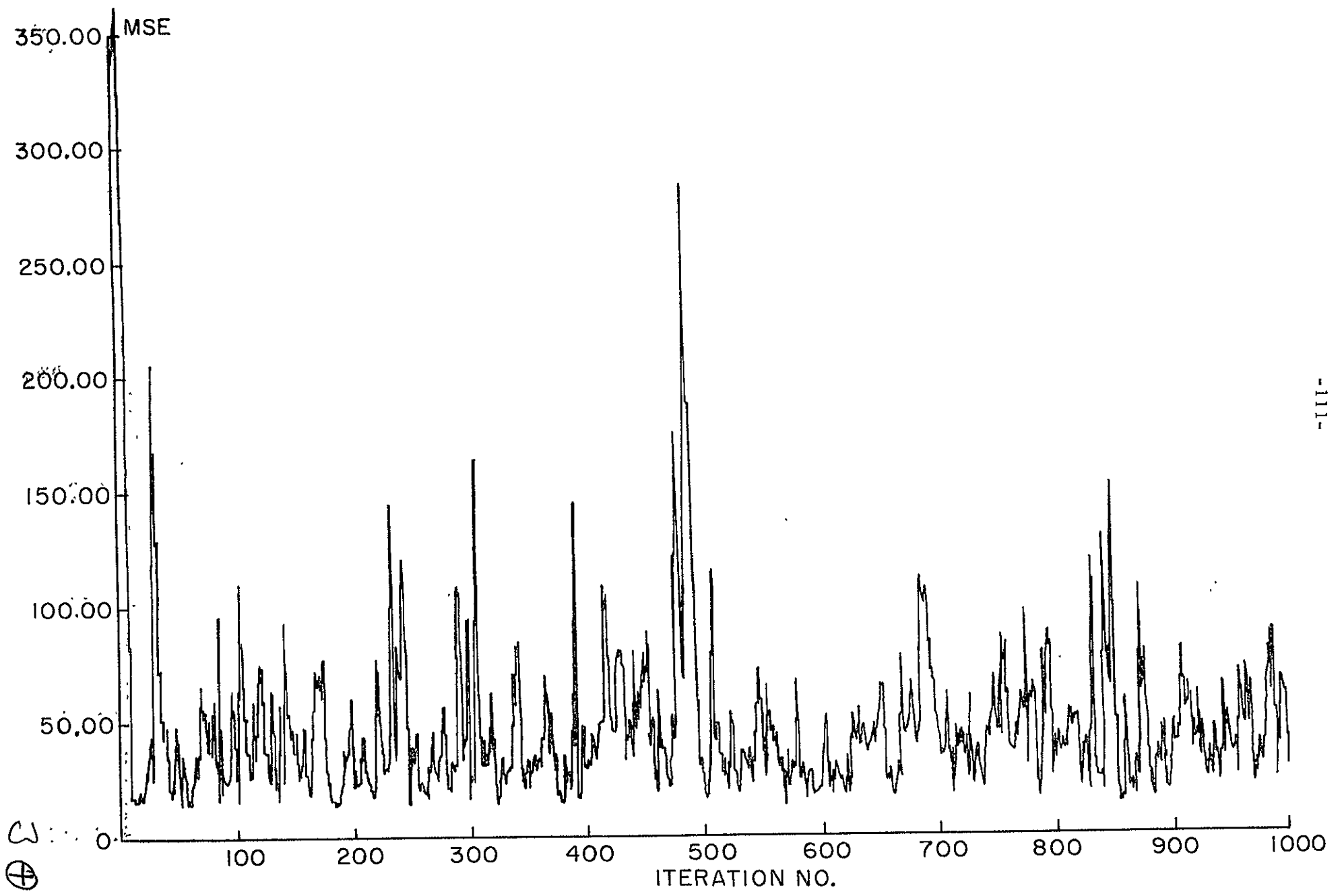


Fig. 4.4.9 Gradient Estimated, Plus Additive Noise

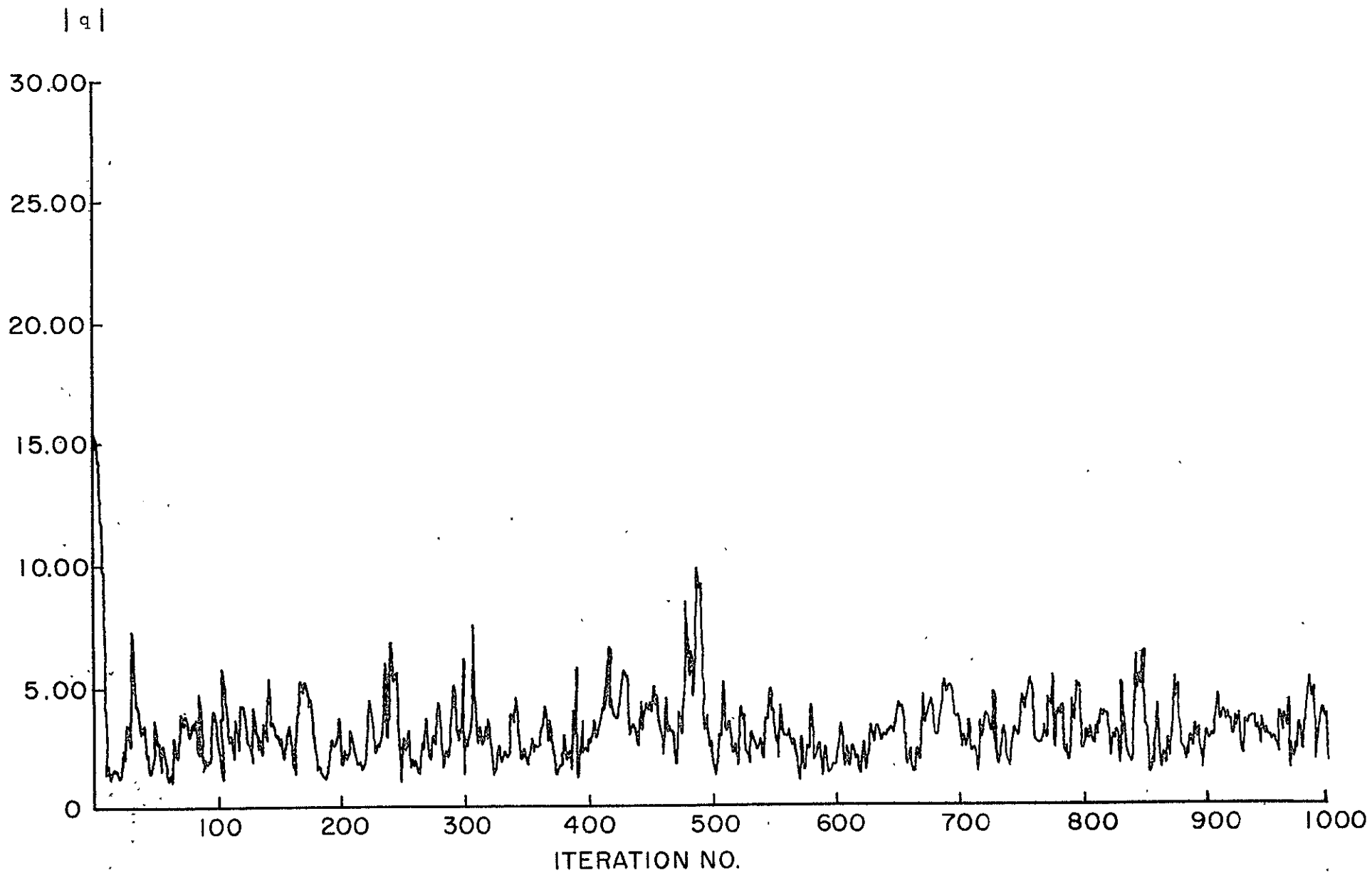


Fig. 4.4.10 Gradient Estimated, Plus Additive Noise



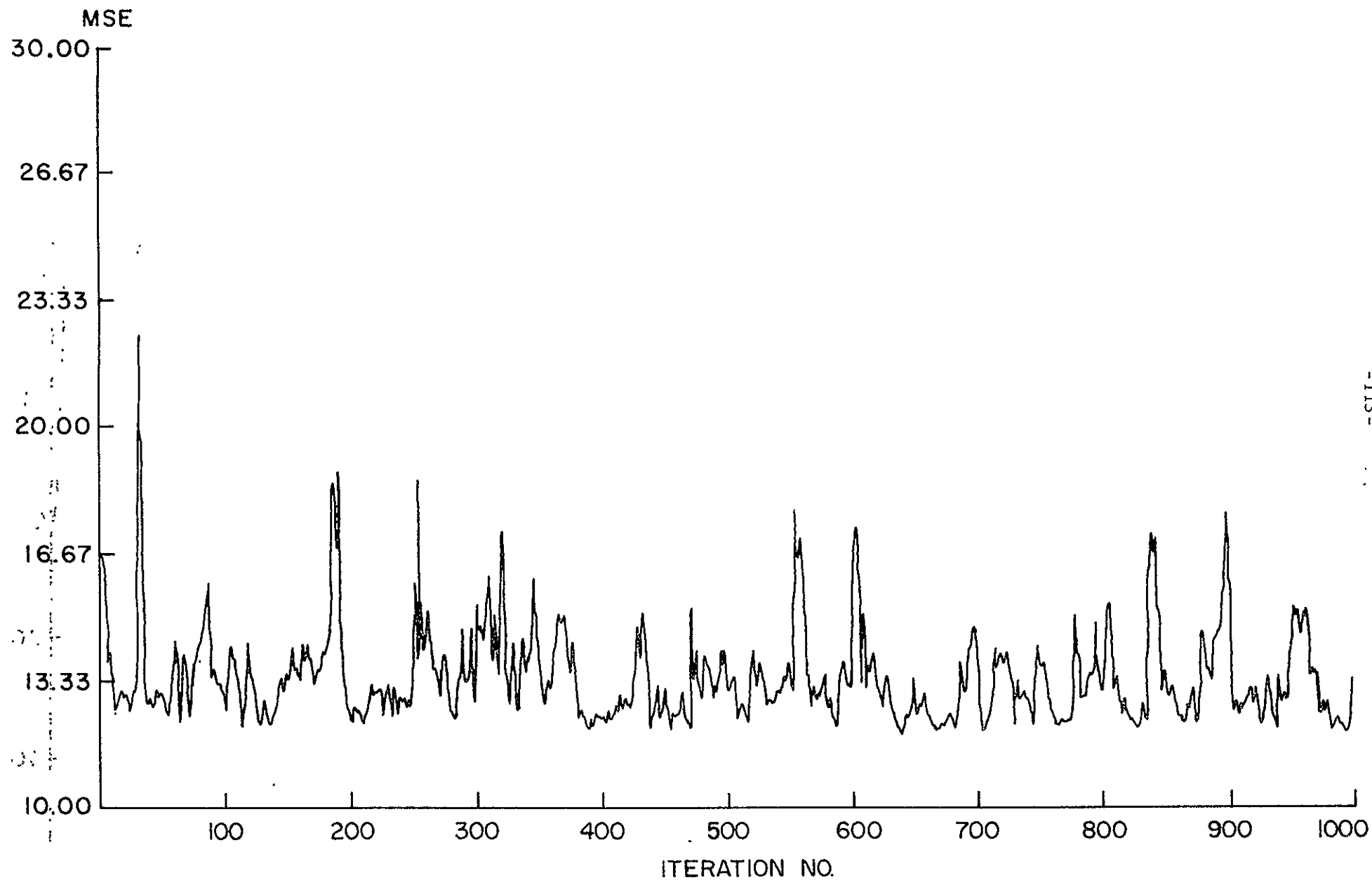


Fig. 4.4.11 Gradient Estimated, Plus Additive Noise

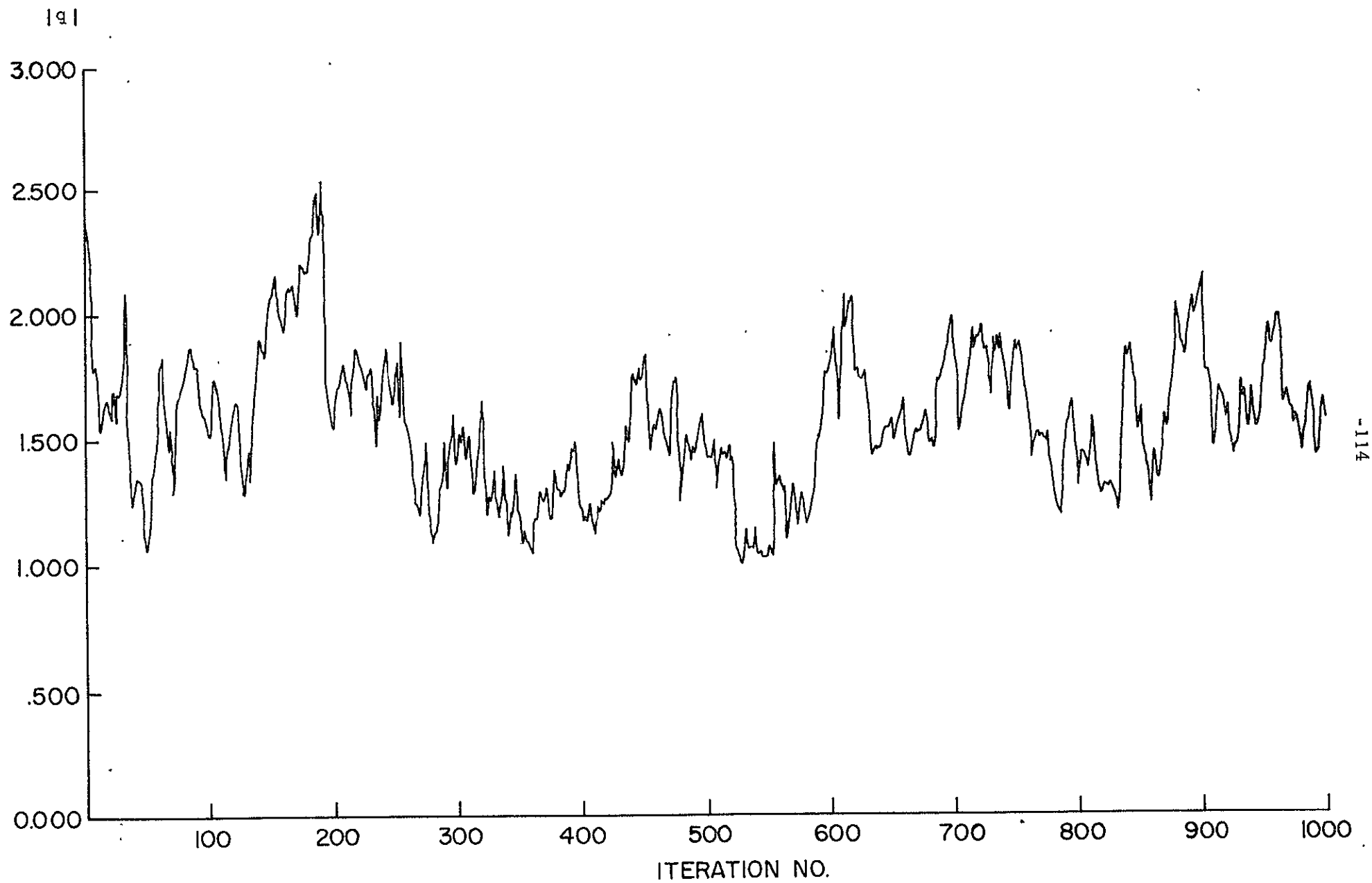


Fig. 4.4.12 Gradient Estimated, Plus Additive Noise

# Appendix A      Proof of Convergence and Bounds on the Asymptotic Variance.

∴ This theorem is essentially the same as Appendix C of Gersho's<sup>(18)</sup> paper.

Theorem: Let  $H_k$  be a sequence of random  $N \times N$  matrices and  $\underline{h}_k$  a sequence of random  $N$ -tuple vectors. Suppose  $E\{H_k\}$  and  $E\{\underline{h}_k\}$  are independent of  $k$ ;  $H_k$  and  $\underline{h}_k$  are independent of  $H_j$  and  $\underline{h}_j$  for  $k \neq j$ ;  $E\{\underline{h}_k\} = \underline{0}$ ; the elements of  $H_k$  and  $\underline{h}_k$  have finite variance;  $E\{H_j\} \equiv \underline{a}$ ,  $\xi_j \equiv ||I - k \underline{a}|| = 1 - k c$  where  $c > 0$ .

Define the random sequence  $\underline{q}_j$  by:

$$\underline{q}_{j+1} = \underline{q}_j - k \underline{\varphi}_j \tag{A 1}$$

$$\underline{\varphi}_j = H_j \underline{q}_j + \underline{h}_j \tag{A 2}$$

for  $j = 0, 1, 2, \dots$  and  $\underline{q}_0$  is an arbitrary deterministic vector. Then for  $k$  positive and sufficiently small

$$\lim_{j \rightarrow \infty} ||E\{\underline{q}_j\}|| = 0 \tag{A 3}$$

and

$$\lim_{j \rightarrow \infty} \sup ||\underline{q}_j|| \leq V(k) \tag{A 4}$$

with  $V(k)$  satisfying

$$\lim_{k \rightarrow 0} V(k) = 0 \tag{A 5}$$

Note that the norm of a random vector  $\underline{u}$  is defined as

$$||\underline{u}|| \equiv \sqrt{E\{\underline{u}^T \underline{u}\}} \tag{A 6}$$

Proof: Combining equations (A1) and (A2) yields

$$\underline{q}_{j+1} = (I - kH_j) \underline{q}_j + k\underline{h}_j \quad (A 7)$$

Since  $\underline{q}_j$  is independent of  $H_j$ , taking the expected value of equation (A7) gives

$$E \{ \underline{q}_{j+1} \} = (I - kA) E \{ \underline{q}_j \} \quad (A 8)$$

Thus

$$\| E \{ \underline{q}_j \} \| \leq \xi^j \| E \{ \underline{q}_0 \} \| \quad (A 9)$$

Since  $\xi < 1$  by hypothesis, equation (A3) follows.

To prove equation (A4), observe that

$$\begin{aligned} E \{ \underline{q}_{j+1}^T \underline{q}_{j+1} \} &= E \{ \underline{q}_j^T (I - kH_j^T) (I - kH_j) \underline{q}_j \} - E \{ \underline{q}_j^T (I - kH_j^T) k\underline{h}_j \} \\ &\quad - E \{ k\underline{h}_j^T (I - kH_j) \underline{q}_j \} + k^2 E \{ \underline{h}_j^T \underline{h}_j \} \end{aligned} \quad (A10)$$

But since  $\underline{q}_j$  is independent of  $H_j$ , the first term in equation (A10) may be bounded by

$$\begin{aligned} E \{ \underline{q}_j^T (I - kH_j^T) (I - kH_j) \underline{q}_j \} &= E \{ \underline{q}_j^T E \{ (I - kH_j^T) (I - kH_j) \} \underline{q}_j \} \\ &\leq \| E \{ (I - kH_j^T) (I - kH_j) \} \| \| \underline{q}_j \|^2 = \mu \| \underline{q}_j \|^2 \end{aligned}$$

$$\text{where } \mu \equiv \| E \{ (I - kH_j^T) (I - kH_j) \} \| \quad (A12)$$

$$\text{Note } \underline{x}^T A \underline{x} \leq \| A \| \| \underline{x} \|^2$$

Combining the second and third terms and using the Schwarz inequality gives

$$\begin{aligned}
 & - 2 k E \{ \underline{q}_j^T (I - k H_j) \underline{h}_j \} \\
 & = - 2 k \left[ E \{ \underline{q}_j^T \} E \{ \underline{h}_j \} - k E \{ \underline{q}_j^T \} E \{ H_j \underline{h}_j \} \right] = 2 k^2 E \{ \underline{q}_j^T \} E \{ H_j \underline{h}_j \} \\
 & \leq 2 k^2 f \| E \{ \underline{q}_j \} \| \quad (A13)
 \end{aligned}$$

$$\text{where} \quad f \equiv \| E \{ H_j \underline{h}_j \} \| \quad (A14)$$

and  $f$  is finite.

Using (A9) we get

$$- 2 k E \{ \underline{q}_j^T (I - k H_j) \underline{h}_j \} \leq 2 k^2 f \xi^j \| E \{ \underline{q}_0 \} \| \quad (A15)$$

Applying the bounds (A11) and (A15) to (A10) yields

$$\| \underline{q}_{j+1} \|^2 = E \{ \underline{q}_{j+1}^T \underline{q}_{j+1} \} \leq \mu \| \underline{q}_j \|^2 + 2 k^2 f \xi^j \| E \{ \underline{q}_0 \} \| + k^2 \| \underline{h}_j \|^2 \quad (A16)$$

If we now define the bounding sequence of positive numbers  $Q_k$  according to

$$Q_0 = \| E \{ \underline{q}_0 \} \|^2 \quad (A17)$$

and

$$Q_{j+1} = \mu Q_j + 2 k^2 f \xi^j \| E \{ \underline{q}_0 \} \| + k^2 \| \underline{h}_j \|^2 \quad (A18)$$

then it follows from (A16) that

$$\| \underline{q}_{j+1} \|^2 \leq Q_k \quad (A19)$$

But the difference equation (A18) has the asymptotic solution.

$$\lim_{j \rightarrow \infty} Q_j = \frac{k^2 \|\underline{h}_j\|^2}{1 - \mu} \quad (\text{A20})$$

because  $\xi < 1$ .

Thus

$$\lim_{j \rightarrow \infty} \sup \|\underline{q}_j\|^2 \leq \frac{k^2 \|\underline{h}_j\|^2}{1 - \mu} \quad (\text{A21})$$

where  $\|\underline{h}_j\|$  is independent of  $j$  by hypothesis.

Let us investigate the positive constant  $\mu$ :

$$\text{if } G_j \equiv H_j - a \quad (\text{A22})$$

$$\text{then } (I - k H_j^T) (I - k H_j) = (I - k a^T - k G_j^T) (I - k a - k G_j)$$

$$= (I - k a^T) (I - k a) - (I - k a^T) k G_j - k G_j^T (I - k a) + k^2 G_j^T G_j$$

$$E \{ (I - k H_j^T) (I - k H_j) \} = E \{ (I - k a^T) (I - k a) \} + k^2 E \{ G_j^T G_j \}$$

$$E \{ (I - k H_j^T) (I - k H_j) \} = (I - k a^T) (I - k a) + k^2 E \{ G_j^T G_j \} \quad (\text{A23})$$

$$\mu = \| (I - k a^T) (I - k a) + k^2 E \{ G_j^T G_j \} \|$$

$$\leq \xi^2 + k^2 \gamma \quad (\text{A24})$$

where

$$\gamma \equiv \| G_j \|^2 \text{ is finite.}$$

Futhermore, in all cases  $\xi$  is of the form  $\xi = 1 - k c$  where  $c > 0$ .

$$\begin{aligned}
 \frac{k^2}{1-\mu} &\leq \frac{k^2}{1 - [(1-kc)^2 + k^2\gamma]} \\
 &= \frac{k^2}{2kc - k^2c^2 - k^2\gamma} = \frac{k}{2c - k(c^2 + \gamma)} \quad (A25)
 \end{aligned}$$

Equations (A4) and (A5) are satisfied if we define

$$V(k) \equiv \frac{k}{2c - k(c^2 + \gamma)} \quad (A26)$$

Q E D

## Appendix B   Rosen's<sup>(21)-(22)</sup> Gradient Projection Algorithm

In this investigation, we indicated that our gradient projection algorithm which adaptively adjusted the tap gains could be thought of as a modification of Rosen's algorithm. Therefore, let us now summarize some well-known linear (and nonlinear) programming methods of optimizing functions subject to linear (and nonlinear) constraints when no noise is present; explain why Rosen's method is applicable to the problem of optimizing functions subject to both constraints and noise; and illustrate, for those unfamiliar with Rosen's algorithm, how it would be used to locate the maximum of a concave function subject to linear constraints.

We restrict our discussion to gradient methods of linear and nonlinear programming because other methods of optimizing convex functions (e.g. Simplex) work essentially by examining the vertices of the feasible region, and testing whether or not the conditions for optimality are satisfied at the vertex being tested. If the conditions are not satisfied we jump to the next vertex. However, since the vertices may be far away from one another, jumping from one vertex to another is not what we want in an adaptive algorithm, which must have the property that if we are not at the exact optimum we must still be "close to" the exact optimum, not at the next vertex which may be a considerable distance away. Another point to consider is that at any single iteration you don't want to move too great a distance because we will sometimes be moving in the wrong direction due to the presence of noise. This is another reason why we don't want to consider just vertices, but rather all points on the boundary of the feasible domain.

All gradient procedures work by moving from an iteration point  $\underline{x}^k$  in the direction of the gradient or, if this is not possible because of the constraints, in the direction of a vector  $\underline{s}$  which makes an acute angle with the gradient, i.e.  $\underline{s}^T \nabla F(\underline{x}^k) > 0$ . We move in this direction until either  $F$  reaches its maximum in this direction or until we cannot go further without leaving the feasible domain. The end point gives the next iteration value  $\underline{x}^{k+1}$ . We never leave the feasible domain throughout the entire iteration.



Zoutendijk's<sup>(24)</sup> method chooses  $\underline{s}$  so that, after a suitable normalization, its scalar product with the gradient is maximized under the condition that we do not immediately leave the feasible domain when moving from  $\underline{x}^k$  in the direction  $\underline{s}$ . We will not use this algorithm because the maximization step uses the abovementioned linear programming methods which are adversely affected by noise. Another procedure is to restrict the vector  $\underline{s}$  to lie in a certain linear manifold of dimension smaller than  $n$ . This approach is used by Rosen. These two methods are somewhat similar. We will use Rosen's method because the iteration steps appear to be simpler and should use less computer time.

We will abstract pp 163-170 from Kunzi, Krelle, and Oettli<sup>(25)</sup> and some numerical examples from Hadley.<sup>(26)</sup> For more details and proofs as well as a discussion of how the algorithm may be modified to account for nonlinear constraints, see Rosen's original papers.

The problem is to maximize the concave function  $F(\underline{x})$  subject to the linear constraints (nonlinear constraints are discussed in Rosen's second paper).

$$h_j(\underline{x}) \equiv \underline{a}_j^T \cdot \underline{x} - b_j \leq 0 \quad j = 1, 2, \dots, m \quad (B\ 1)$$

where  $\underline{x}$  is an  $n$  dimensional vector.

If a point  $\underline{x}^0$  of the feasible domain (i. e.  $\underline{x}^0$  satisfies all the constraints) is not the constrained maximum, then we may look for another feasible point with a higher function value by proceeding from  $\underline{x}^0$  in the direction of the gradient of the objective function. This is always possible if  $\underline{x}^0$  is an interior point. However, the method can fail if  $\underline{x}^0$  is a boundary point, because the gradient vector may point toward the exterior of the feasible domain. Rosen's method is to project the gradient onto the boundary of the feasible domain and then proceed in the direction of the projection rather than in the direction of the gradient itself. More precisely, the gradient is projected onto a linear submanifold of the boundary, i. e. on the submanifold of least dimension that contains  $\underline{x}^0$ . In three dimensional space, for instance, the feasible domain is a polyhedron whose boundary consists of manifolds of dimension two (faces), dimension one (edges), and dimension

zero (vertices). If  $\underline{x}^0$  lies on a face but not on an edge, the gradient is projected onto this face; if  $\underline{x}^0$  lies on an edge, we project on the edge. Rosen's method coincides with the usual gradient method if the point  $\underline{x}^0$  lies in the interior of the feasible domain.

We denote the  $(n-1)$  dimensional manifold (boundary hyperplane) defined by  $h_j(\underline{x}) = 0$  by  $H_j$ , i. e.

$$H_j \equiv \{ \underline{x} \mid h_j(\underline{x}) = 0 \} \quad j = 1, 2, \dots, m \quad (B \ 2)$$

The boundary of the feasible domain consists of all feasible points  $[\underline{x} \mid h_j(\underline{x}) \leq 0 \text{ for all } j]$  with  $h_j(\underline{x}) = 0$  for at least one  $j$ . The (non-normalized) normal vector  $\underline{a}_j$  is perpendicular to  $H_j$  and points outward from the feasible domain. A number of hyperplanes  $H_j$  are linearly independent if the corresponding  $\underline{a}_j$  are linearly independent. The intersection of  $k$  hyperplanes is the set of points which lie simultaneously on all  $k$  hyperplanes. The intersection of  $k$  linearly independent hyperplanes forms an  $(n-k)$  dimensional linear manifold in the  $n$  dimensional space of the  $\underline{x}$  vectors.

Let us now consider the projection of the gradient vector. Say  $\underline{x}^0$  lies on  $r$  hyperplanes. We pick out  $q$  linearly independent hyperplanes from among these  $r$ , which, after a suitable reordering of the indices we may assume to be  $H_1, \dots, H_q$ . Let  $D$  denote the  $(n-q)$  dimensional intersection of these hyperplanes. The normals  $\underline{a}_1, \dots, \underline{a}_q$  are perpendicular to the linear manifold  $D$ . The  $q$  dimensional linear manifold spanned by  $\underline{a}_1, \dots, \underline{a}_q$  will be denoted by  $\tilde{D}$ .  $D$  and  $\tilde{D}$  are mutually perpendicular and together span the whole space. The projection of a vector  $\underline{y}$  on the linear manifold  $D$  is denoted by  $\underline{y}_D$  and is given by

$$\underline{y}_D \equiv P_q \underline{y} \quad (B \ 3)$$

$$\text{where} \quad P_q \equiv I - A_q (A_q^T A_q)^{-1} A_q^T \quad (B \ 4)$$

$$\text{and} \quad A_q \equiv (\underline{a}_1 \ \underline{a}_2 \ \dots \ \underline{a}_q) \quad (B \ 5)$$

Note that  $P_0 \equiv I$  and  $P_n \equiv$  zero matrix.

Rosen proves that the point  $\underline{x}^k$  is the unique constrained maximum for concave objective functions if and only if  $\underline{x}^k$  satisfies

$$P_q \underline{g}(\underline{x}^k) = \underline{0} \quad (B\ 6)$$

and

$$(A_q^T A_q)^{-1} A_q^T \underline{g}(\underline{x}^k) \geq \underline{0} \quad (B\ 7)$$

where  $\underline{g}(\underline{x}^k)$  is the gradient vector at point  $\underline{x}^k$ .

Condition (B6) states that the gradient vector is orthogonal to the manifold  $D$ , and thus lies in  $\tilde{D}$ . Hence

$$\underline{g}(\underline{x}^k) = \sum_{j=1}^q u_j \underline{a}_j = A_q \underline{u} \quad (B\ 8)$$

Substituting (B8) into (B7) we see that (B7) may be rewritten as

$$\underline{u} \geq \underline{0}$$

Equations (B6) and (B7) together imply that a necessary and sufficient condition for the point  $\underline{x}^k$  to be a constrained maximum is that the gradient of the objective function be expressible as a non-negative linear combination of the exterior normals to the hyperplanes on which the point lies. This is equivalent to the well-known Kuhn-Tucker<sup>(27)</sup> conditions. If  $\underline{x}^k$  is an interior point of the feasible domain, the optimality criterion simplifies to  $P_0 \underline{g}(\underline{x}^k) = \underline{g}(\underline{x}^k) = \underline{0}$ .

Whenever the conditions for optimality are not satisfied Rosen shows there exists a feasible point  $\underline{x}^{k+1}$  which yields a higher objective function value. There are two possibilities (we avoid discussing degeneracies) which we consider separately. Denote  $\underline{g}(\underline{x}^k)$  by  $\underline{g}_k$ .

Case I  $P_q \underline{g}_k \neq \underline{0}$ .

This means that  $\underline{x}^0$  is not a vertex of the feasible domain, i.e.  $q < n$ , and  $D$  has at least the dimension of a straight line. We move in

the direction given by the vector  $\underline{s}^k = P_q \underline{g}_k$  (B9). We will not discuss here how far to move in this direction because this part of Rosen's algorithm does not apply to our modification of Rosen's algorithm.

$$\text{Case II} \quad P_q \underline{g}_k = \underline{0}$$

but  $u_j < 0$  for at least one  $j$ . We then choose one of the indices for which  $u_j < 0$ , e.g. the one for which  $|\underline{a}_j| u_j$  is most negative, and then disregard the corresponding hyperplane  $H_j$ . Suppose this is the hyperplane  $H_q$ . Then  $u_q < 0$ , and we proceed as if  $\underline{x}^k$  lies only on  $H_1$  to  $H_{q-1}$ , i.e. we raise the dimension of  $D$  by one. The associated projection matrix is now  $P_{q-1}$ . We have  $P_{q-1} \underline{a}_q \neq \underline{0}$  because  $\underline{a}_q$  is independent of  $\underline{a}_1$  to  $\underline{a}_{q-1}$ . This implies that

$$P_{q-1} \underline{g}_k = P_{q-1} \left( \underline{z} + \sum_{j=1}^q u_j \underline{a}_j \right) = u_q P_{q-1} \underline{a}_q \neq \underline{0}$$

where  $\underline{z}$  belongs to  $D$ . Consequently, in the new  $D$ , which has one dimension more, we have the same situation as in case I, and we can proceed as in that case by setting

$$\underline{s}^k = P_{q-1} \underline{g}_k \quad (\text{B10})$$

These are the main steps involved in Rosen's algorithm. We add that nonlinear constraints can also be handled, but we will not discuss that algorithm (see Rosen's papers, and chapter six of this investigation) here.

Finally we present two examples, taken from Hadley, to illustrate how the algorithm works. Consider Fig B1

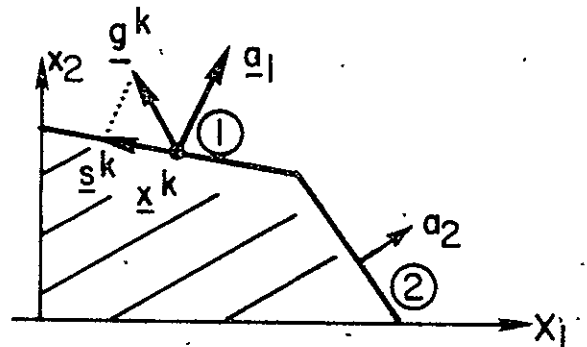


Fig. B1 Diagram for example one

Assume that the current feasible solution is  $\underline{x}^k$ . We cannot move in the direction of the gradient without violating constraint 1. The vector  $\underline{s}^k$  is given by (B9)

$$\begin{aligned}\underline{s}^k &= P_1 \underline{g}^k = \left[ I - \frac{\underline{a}_1 \underline{a}_1^T}{|\underline{a}_1|^2} \right] \underline{g}^k \\ &= \underline{g}^k - \left( \frac{\underline{a}_1^T \underline{g}^k}{|\underline{a}_1|^2} \right) \underline{a}_1\end{aligned}$$

This is nothing more than the perpendicular projection of  $\underline{g}^k$  onto the boundary of the set of feasible solutions, as shown.

Consider next the situation illustrated in Fig B2

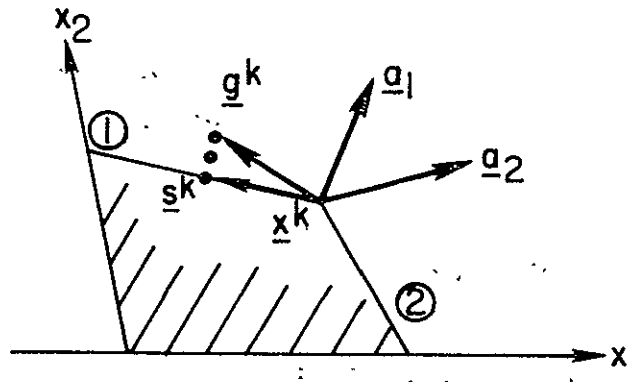


Fig. B2 Diagram for example two

Both constraints will be violated if we move in the direction of the gradient vector. Also  $P_2 \underline{g}^k = \underline{0}$  indicating that it is not possible to move from  $\underline{x}^k$  in any direction such that both constraints hold as strict equalities. Note that when  $\underline{g}^k$  is expressed as a linear combination of  $\underline{a}_1$  and  $\underline{a}_2$ ,  $\underline{g}^k = \alpha_1 \underline{a}_1 + \alpha_2 \underline{a}_2$  we see that  $\alpha_2$  is negative. We can find a feasible direction in which to move (case II) by allowing constraint 2 to hold as a strict inequality, while constraint 1 holds as a strict equality. If we do this, the problem is reduced to the previous illustration.

## CHAPTER 5

### Soft Constraints

#### Section 5.1 Introduction

In the last chapter, we devised an algorithm that minimizes an objective function subject to constraints which were never to be violated. In this chapter, we will devise an algorithm that differs from the gradient projection algorithm of the previous chapter in that this algorithm minimizes an objective function subject to constraints which may be "slightly" violated, but which cannot be violated "too much." This type of constraint is known in the literature as a "soft" constraint as opposed to the "hard" constraint dealt with in chapter four.

Again, our final objective is to design an adaptive algorithm which will maximize the SNR subject to a constraint on the super-gain ratio when unknown interfering noise is present. Again because the SNR and super-gain ratios are nonlinear quantities, it is difficult to prove convergence of our algorithm or to analytically find the algorithm's rate of convergence. Again, for the purpose of mathematical tractability and because it is useful in its own right, we will consider an adaptive algorithm which minimizes the MSE subject to a linear constraint.

The algorithms of this chapter are simply a gradient minimization of a convex modified objective function, the modified objective function consisting of our original objective function plus a convex penalty function which serves to increase the value of our modified objective function whenever the constraints are violated, i. e. we will minimize the convex function

$$f(\underline{W}) = \epsilon_{\min}^2 + (\underline{W}^T - \underline{W}_{LMS}^T) \phi(\underline{W} - \underline{W}_{LMS}) \quad (5.1.1)$$

subject to the "soft" linear constraint, shown in Fig 5.1.1 below

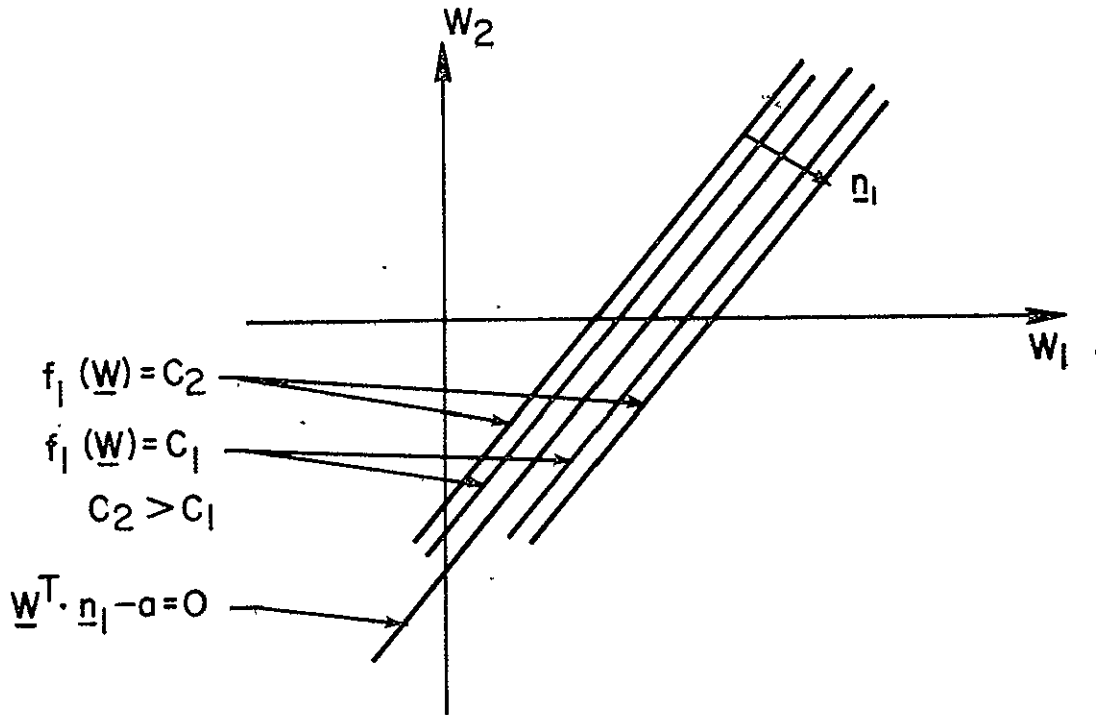


Fig. 5.1.1 Constraint and Penalty Function Level Curves

The constraint equation is of the form

$$\underline{W}^T \cdot \underline{n}_1 - a = 0 \quad (5.1.2)$$

The convex penalty function we will use is given by

$$f_1(\underline{W}) = K_1 [\underline{W}^T \cdot \underline{n}_1 - a]^2 \quad (5.1.3)$$

The level curves of this penalty function are also shown in Fig 5.1.1.

We should note that if  $K_1$  is "large enough" we will always be very "close" to the line  $\underline{W}^T \cdot \underline{n}_1 - a = 0$  which then may be interpreted as a linear approximation (i.e. the first terms of a Taylor expansion) at point  $\underline{W}$  to any arbitrary nonlinear constraint (e.g. the super-gain ratio) provided that as the algorithm moves from point to point in the  $\underline{W}$  space, we keep replacing the nonlinear constraint by the best linear approximation to it at each point.

Assuming we have only one constraint in the problem, as given by equation (5.1.2) we will present three algorithms, corresponding to the

three cases studied in chapter four, i.e. when the gradient is known, when we have a noise-free estimate of the gradient, and when we have a noisy estimate of the gradient, and for each of these algorithms we will investigate convergence (convergence of the expected value of the weight vectors and bounds on the variance of the weight vectors in cases two and three), the rate of convergence, and the bias between what our "soft" constraint algorithms converge to and the optimum weight vector when we have a "hard" constraint, which was found in section 4.2 to be

$$\underline{W}_{opt} = \underline{W}_{LMS} + \frac{(\underline{a} - \underline{W}_{LMS}^T \cdot \underline{n}_1)}{(\underline{n}_1^T \phi^{-1} \underline{n}_1)} \phi^{-1} \underline{n}_1 \quad (5.1.4)$$

All three algorithms seek to minimize the modified convex objective ( $j$  indicates the iteration number)

$$F(\underline{W}_j) = \epsilon_{min}^2 + (\underline{W}_j^T - \underline{W}_{LMS}^T) \phi (\underline{W}_j - \underline{W}_{LMS}) + K_1 [\underline{W}_j^T \cdot \underline{n}_1 - a]^2 \quad (5.1.5)$$

In case 1, the gradient of equation (5.1.5) is

$$\underline{g}(\underline{W}_j) = 2 \phi (\underline{W}_j - \underline{W}_{LMS}) + 2 K_1 [\underline{W}_j^T \cdot \underline{n}_1 - a] \underline{n}_1 \quad (5.1.6)$$

In case 2, we assume  $\phi$  is not available and must be estimated by  $\underline{s}_j$ ,  $\underline{d}_j$  and  $\underline{W}_j$  which are available

$$\underline{g}(\underline{W}_j) = -2 \underline{s}_j (\underline{d}_j - \underline{s}_j^T \cdot \underline{W}_j) + 2 K_1 [\underline{W}_j^T \cdot \underline{n}_1 - a] \underline{n}_1 \quad (5.1.7)$$

In case 3, we assume  $\underline{s}_j$  is not available, but a noisy estimate of  $\underline{s}_j$  is available

$$\underline{g}(\underline{W}_j) = -2(\underline{s}_j + \underline{n}_j) [\underline{d}_j - (\underline{s}_j^T + \underline{n}_j^T) \underline{W}_j] + 2 K_1 [\underline{W}_j^T \cdot \underline{n}_1 - a] \underline{n}_1 \quad (5.1.8)$$



Section 5.2.1      The Algorithm, Proof of Convergence, and Bounds  
on the Rate of Convergence if the Gradient is Known.

Using equation (5.1.6) the algorithm is

$$\underline{W}_{j+1} = \underline{W}_j - k \left\{ 2\phi(\underline{W}_j - \underline{W}_{LMS}) + 2K_1 [\underline{W}_j^T \cdot \underline{n}_1 - a] \underline{n}_1 \right\} \quad (5.2.1.1)$$

The above equations are a set of first order deterministic difference equations. Let us first solve for the asymptotic value of  $\underline{W}$ , denoted by  $\underline{W}_\infty$ . Setting  $\underline{W}_{j+1} = \underline{W}_j = \underline{W}_\infty$  gives

$$\underline{W}_\infty = \underline{W}_{LMS} - K_1 [\underline{W}_\infty^T \cdot \underline{n}_1 - a] \phi^{-1} \underline{n}_1 \quad (5.2.1.2)$$

$$\text{Let } \underline{W}_\infty = \underline{c} + d \cdot \underline{n}_1 \quad (5.2.1.3)$$

$$\text{where } \underline{c}^T \cdot \underline{n}_1 = \underline{n}_1^T \cdot \underline{c} = 0 \quad (5.2.1.4)$$

Remembering that  $\underline{n}_1^T \cdot \underline{n}_1 = 1$  we have

$$\underline{c} + d \underline{n}_1 = \underline{W}_{LMS} - K_1 [d - a] \phi^{-1} \underline{n}_1 \quad (5.2.1.5)$$

Multiplying by  $\underline{n}_1^T$  on the left yields

$$d = \frac{1}{1 + K_1 (\underline{n}_1^T \phi^{-1} \underline{n}_1)} \left[ \underline{n}_1^T \cdot \underline{W}_{LMS} + K_1 a (\underline{n}_1^T \phi^{-1} \underline{n}_1) \right] \quad (5.2.1.6)$$

Substituting (5.2.1.6) into (5.2.1.5) yields

$$\underline{c} = -d \underline{n}_1 + \underline{W}_{LMS} - K_1 \left[ \frac{\underline{n}_1^T \cdot \underline{W}_{LMS} + K_1 a (\underline{n}_1^T \phi^{-1} \underline{n}_1)}{1 + K_1 (\underline{n}_1^T \phi^{-1} \underline{n}_1)} - a \right] \phi^{-1} \underline{n}_1$$

and

$$\underline{W}_{\infty} = \underline{W}_{LMS} + \frac{K_1}{1 + K_1 (\underline{n}_1^T \phi^{-1} \underline{n}_1)} \left[ a - \underline{n}_1^T \cdot \underline{W}_{LMS} \right] \phi^{-1} \underline{n}_1 \quad (5.2.1.7)$$

If we let  $K_1 \rightarrow \infty$ , which means that the penalty function is infinite unless the weight vector lies exactly on the line  $\underline{W}^T \underline{n}_1 - a_1 = 0$ ,  $\underline{W}_{\infty}$  becomes

$$\underline{W}_{\infty} = \underline{W}_{LMS} + \frac{1}{(\underline{n}_1^T \phi^{-1} \underline{n}_1)} \left[ a - \underline{n}_1^T \cdot \underline{W}_{LMS} \right] \phi^{-1} \underline{n}_1$$

which is the optimum solution in the "hard" constraint case (see equation (5.1.4)).

By comparing equation (5.2.1.7) which tells us the steady state value of  $\underline{W}$  that our algorithm converges to, and equation (5.1.4) which tells us what the optimum value that we want to converge to is, we can get an idea of how to choose  $K_1$ , i.e. in the steady state our penalty algorithm converges to  $\underline{W}_{\infty} = \underline{W}_{LMS} + \underline{x}$  where the direction of the vector  $\underline{x}$  is the same as the direction of  $\underline{x}_{opt}$  where  $\underline{W}_{opt} = \underline{W}_{LMS} + \underline{x}_{opt}$ ; however the magnitude of  $\underline{x}$  is less than the magnitude of  $\underline{x}_{opt}$ . If we want this bias to be less than, say 1%, we must choose  $K_1$  to satisfy

$$\frac{K_1}{1 + K_1 (\underline{n}_1^T \phi^{-1} \underline{n}_1)} = \frac{.99}{(\underline{n}_1^T \phi^{-1} \underline{n}_1)}$$

which implies

$$K_1 \geq \frac{99}{(\underline{n}_1^T \phi^{-1} \underline{n}_1)} \geq \frac{99}{\rho_1}$$

where  $\rho_1$  is the minimum eigenvalue of  $\phi^{-1}$ .

We will now investigate how fast our algorithm converges to  $\underline{W}_\infty$ .

$$\text{Define } \underline{q}_j \equiv \underline{W}_j - \underline{W}_\infty \quad (5.2.1.8)$$

In terms of  $\underline{q}$ , the algorithm is

$$\underline{q}_{j+1} = \underline{q}_j - k \left\{ 2 \phi(\underline{q}_j + \underline{W}_\infty - \underline{W}_{\text{LMS}}) + 2K_1 \left[ \underline{q}_j^T \cdot \underline{n}_1 + \underline{W}_\infty^T \cdot \underline{n}_1 - a \right] \underline{n}_1 \right\} \quad (5.2.1.9)$$

$$\begin{aligned} \underline{q}_{j+1} = \underline{q}_j - 2k \phi \underline{q}_j - 2k K_1 (\underline{q}_j^T \cdot \underline{n}_1) \underline{n}_1 \\ - 2k \phi(\underline{W}_\infty - \underline{W}_{\text{LMS}}) - 2k K_1 \underline{n}_1 (\underline{n}_1^T \cdot \underline{W}_\infty - a) \end{aligned} \quad (5.2.1.10)$$

$$\text{But } \underline{W}_\infty - \underline{W}_{\text{LMS}} = \frac{K_1}{1 + K_1 (\underline{n}_1^T \phi^{-1} \underline{n}_1)} \left[ a - \underline{n}_1^T \cdot \underline{W}_{\text{LMS}} \right] \phi^{-1} \underline{n}_1$$

$$- 2k \phi(\underline{W}_\infty - \underline{W}_{\text{LMS}}) = \frac{- 2k K_1}{1 + K_1 (\underline{n}_1^T \phi^{-1} \underline{n}_1)} \left[ a - \underline{n}_1^T \cdot \underline{W}_{\text{LMS}} \right] \underline{n}_1$$

and

$$\begin{aligned} \underline{n}_1^T \cdot \underline{W}_\infty &= \underline{n}_1^T \cdot \underline{W}_{\text{LMS}} + \frac{K_1 (a - \underline{n}_1^T \cdot \underline{W}_{\text{LMS}}) (\underline{n}_1^T \phi^{-1} \underline{n}_1)}{1 + K_1 (\underline{n}_1^T \phi^{-1} \underline{n}_1)} \\ &= \frac{\underline{n}_1^T \cdot \underline{W}_{\text{LMS}} + K_1 (\underline{n}_1^T \phi^{-1} \underline{n}_1) a}{1 + K_1 (\underline{n}_1^T \phi^{-1} \underline{n}_1)} \end{aligned}$$

Equation (5.2.1.9) then becomes

$$\underline{q}_{j+1} = \underline{q}_j - 2k \phi \underline{q}_j - 2k K_1 \underline{n}_1 (\underline{n}_1^T \cdot \underline{q}_j) \quad (5.2.1.11)$$

thus

$$\| \underline{q}_{j+1} \| \leq \xi^{j+1} \| \underline{q}_0 \| \quad (5.2.1.12)$$

where

$$\xi \equiv \| I - 2k(\phi + K_1 \underline{n}_1 \underline{n}_1^T) \| \quad (5.2.1.13)$$

Note that  $\phi + K_1 \underline{n}_1 \underline{n}_1^T$  is positive definite symmetric

$$\text{Pf: } \underline{x}^T (\phi + K_1 \underline{n}_1 \underline{n}_1^T) \underline{x} = \underline{x}^T \phi \underline{x} + K_1 \underbrace{(\underline{x}^T \underline{n}_1)(\underline{n}_1^T \underline{x})}_{|\underline{x}^T \underline{n}_1|^2}$$

From Goldstein<sup>(28)</sup> page 24

$$\xi = \max \left\{ |1 - 2k\rho_n|, |1 - 2k\rho_1| \right\} \quad (5.2.1.14)$$

where  $\rho_1$  and  $\rho_n$  are the min and max eigenvalues of  $(\phi + K_1 \underline{n}_1 \underline{n}_1^T)$  respectively. For  $k$  small enough

$$0 < \xi < 1 \quad (5.2.1.15)$$

Equation (5.2.1.12) shows that the rate of convergence is given by the number  $\xi$ , which for  $k$  small enough is between zero and one (thus guaranteeing convergence), and  $\xi \rightarrow 1$  as  $k \rightarrow 0$  (i.e. the rate of convergence becomes slower as  $k \rightarrow 0$ ).

In this section we have proven that our algorithm converges to  $\underline{W}_\infty$  for  $k$  sufficiently small. In the next section we will investigate a more useful algorithm, i.e. an algorithm that does not require a priori knowledge of  $\phi$ .

Section 5.2.2      The Algorithm, Proof of Convergence, and Bounds  
on the Rate of Convergence of the Gradient is Estimated.

Using equation (5.1.7) the algorithm is

$$\underline{W}_{j+1} = \underline{W}_j - k \left\{ -2 \underline{s}_j (d_j - \underline{s}_j^T \underline{W}_j) + 2 K_1 [\underline{W}_j^T \underline{n}_1 - a] \underline{n}_1 \right\} \quad (5.2.2.1)$$

These equations constitute a set of first-order stochastic difference equations. We will first solve for the asymptotic expected value of  $\underline{W}$ , denoted by  $\underline{W}_\infty$ .

Taking the expected value of equation (5.2.2.1) yields

$$\begin{aligned} E\{\underline{W}_{j+1}\} = E\{\underline{W}_j\} - k \left\{ -2 E\{\underline{s}_j d_j\} + 2 E\{\underline{s}_j \underline{s}_j^T \underline{W}_j\} \right. \\ \left. + 2 K_1 [E\{\underline{W}_j^T\} \cdot \underline{n}_1 - a] \underline{n}_1 \right\} \end{aligned}$$

Noting that  $E\{\underline{s}_j \underline{s}_j^T \underline{W}_j\} = E\{\underline{s}_j \underline{s}_j^T\} E\{\underline{W}_j\} = \phi(\underline{s}, \underline{s}) E\{\underline{W}_j\}$  as in chapter four, we may rewrite this equation as

$$\begin{aligned} E\{\underline{W}_{j+1}\} = E\{\underline{W}_j\} - k \left\{ -2 \phi(\underline{s}, d) + 2 \phi(\underline{s}, \underline{s}) E\{\underline{W}_j\} \right. \\ \left. + 2 K_1 [E\{\underline{W}_j^T\} \cdot \underline{n}_1 - a] \underline{n}_1 \right\} \end{aligned}$$

Using equation (4.1.7a)

$$\phi(\underline{s}, d) = \phi(\underline{s}, \underline{s}) \underline{W}_{LMS}$$

$$E\{\underline{W}_{j+1}\} = E\{\underline{W}_j\} + 2k \phi \left( \underline{W}_{LMS} - E\{\underline{W}_j\} \right) - 2k K_1 [E\{\underline{W}_j^T\} \cdot \underline{n}_1 - a] \underline{n}_1$$

We now have a set of deterministic first-order difference equations whose asymptotic value  $E\{\underline{W}_\infty\} \equiv \underline{\bar{W}}_\infty$ , can be found by setting  $E\{\underline{W}_j\} = E\{\underline{W}_\infty\} = \underline{\bar{W}}_\infty$ , giving

$$\underline{\bar{W}}_\infty - \underline{W}_{LMS} = -K_1 \left[ \begin{array}{c} \underline{\bar{W}}_\infty^T \\ \underline{n}_1^T \end{array} \cdot \underline{n}_1 - a \right] \phi^{-1} \underline{n}_1 \quad (5.2.2.2)$$

This is the same as equation (5.2.1.2) and the solution is given by equation (5.2.1.7)

$$\underline{\bar{W}}_\infty = \underline{W}_{LMS} + \frac{K_1}{1 + K_1 (\underline{n}_1^T \phi^{-1} \underline{n}_1)} \left[ a - \underline{n}_1^T \cdot \underline{W}_{LMS} \right] \phi^{-1} \underline{n}_1 \quad (5.2.2.3)$$

Because our difference equations describing the behavior of the weight vectors are stochastic, the above result is not sufficient to prove convergence of the weight vectors to  $\underline{\bar{W}}_\infty$ , we must also show that the variance of the stochastic vectors  $\underline{q}_j \equiv \underline{W}_j - \underline{\bar{W}}_\infty$  is bounded. To do this

$$\text{Define } \underline{q}_j \equiv \underline{W}_j - \underline{\bar{W}}_\infty \quad (5.2.2.4)$$

the algorithm (5.2.2.1) may be rewritten in terms of  $\underline{q}_j$  as

$$\begin{aligned} \underline{q}_{j+1} &= \underline{q}_j - k \left[ 2 \underline{s}_j \underline{s}_j^T + 2 K_1 \underline{n}_1 \underline{n}_1^T \right] \underline{q}_j \\ &\quad - k \left[ 2 \underline{s}_j \underline{s}_j^T + 2 K_1 \underline{n}_1 \underline{n}_1^T \right] \underline{\bar{W}}_\infty \\ &\quad + k \left[ 2 \underline{s}_j \underline{d}_j + 2 K_1 a \underline{n}_1 \right] \end{aligned}$$

$$\text{Define } H_j \equiv 2 \underline{s}_j \underline{s}_j^T + 2 K_1 \underline{n}_1 \underline{n}_1^T \quad (5.2.2.5)$$

$$\underline{V}_j \equiv 2 \underline{s}_j \underline{d}_j + 2 K_1 a \underline{n}_1 \quad (5.2.2.6)$$

$$\text{Thus } \underline{q}_{j+1} = \underline{q}_j - k \underline{\phi}_j \quad (5.2.2.7)$$

$$\text{where} \quad \underline{\varphi}_j \equiv H_j \underline{q}_j + \underline{h}_j \quad (5.2.2.8)$$

$$\text{and} \quad \underline{h}_j \equiv H_j \bar{W}_\infty - \underline{V}_j \quad (5.2.2.9)$$

Note that  $E\{H_j\}$  and  $E\{\underline{h}_j\}$  are independent of  $j$ . Also  $H_j$  and  $\underline{h}_j$  are statistically independent of  $H_k$  and  $\underline{h}_k$  if  $j \neq k$ , because we assumed that  $\underline{s}_j, \underline{s}_k$  are statistically independent for  $k \neq j$ .

Noting that

$$E\{H_j\} = 2\phi + 2K_1 \underline{n}_1 \underline{n}_1^T$$

$$E\{\underline{V}_j\} = 2\underline{\phi}(\underline{s}, d) + 2K_1 \underline{a} \underline{n}_1 = 2\underline{\phi}_{LMS} + 2k_1 \underline{a} \underline{n}_1$$

we may show that

$$E\{\underline{h}_j\} = \underline{0} \quad (5.2.2.10)$$

Note that  $E\{H_j\} = 2\phi + 2K_1 \underline{n}_1 \underline{n}_1^T \equiv \underline{A}$  is a symmetric positive definite matrix.

The algorithm is thus

$$\underline{q}_{j+1} = \underline{q}_j - k \underline{\varphi}_j$$

where

$$\underline{\varphi}_j = H_j \underline{q}_j + \underline{h}_j$$

and  $H_j$  is a sequence of random  $n \times n$  matrices;  $\underline{h}_j$  is a sequence of random  $n$ -tuple vectors; the expected values of  $H_j$  and  $\underline{h}_j$  were shown to be independent of  $j$ ;  $H_j$  and  $\underline{h}_j$  are independent of  $H_\ell$  and  $\underline{h}_\ell$  for  $j \neq \ell$ ;  $E\{\underline{h}_j\} = \underline{0}$ ; and the elements of  $H_j$  and  $\underline{h}_j$  have finite variance, with  $E\{H_j\} = \underline{A}$ , where  $\underline{A}$  is a symmetric positive definite matrix.

Under these conditions, it is shown in Appendix A of chapter four that for  $k$  sufficiently small

$$\lim_{j \rightarrow \infty} ||E\{\underline{q}_j\}|| = 0 \quad (5.2.2.11)$$

$$\text{and} \quad \lim_{j \rightarrow \infty} \sup || \underline{q}_j || \leq V(k) \quad (5.2.2.12)$$

where the norm of a random vector  $\underline{u}$  is defined as

$$|| \underline{u} || \equiv \sqrt{E \{ \underline{u}^T \underline{u} \}}$$

$$\text{and} \quad \lim_{k \rightarrow 0} V(k) = 0 \quad (5.2.2.13)$$

Equation (5.2.2.11) shows again that the random weight vectors converge, in the mean, to  $\underline{W}_{\text{optimum}}$  and (5.2.2.12) shows that the variance of the random weight vectors about their expected value is bounded, and the bound can be made as small as desired by choosing  $k$  sufficiently small as shown by (5.2.2.13).

The rate of convergence of the mean of the random weight vectors is shown in the proof of the above theorem to be bounded by  $\xi$ , where

$$\xi = || I - k ( 2 \phi + 2 K_1 \underline{n}_1 \underline{n}_1^T ) || \quad (5.2.2.14)$$

Since  $A \equiv ( 2 \phi + 2 K_1 \underline{n}_1 \underline{n}_1^T )$  is positive definite symmetric, we have

$$\xi = \max \left\{ | 1 - k \rho_1 | , | 1 - k \rho_n | \right\} \quad (5.2.2.15)$$

where  $\rho_1$  and  $\rho_n$  are both positive, and represent the minimum and maximum eigenvalues of  $A$  respectively, as shown by Goldstein<sup>(28)</sup> page 24.

Thus  $0 < \xi < 1$ , and this again proves convergence of the algorithm of this section. In the next section we will investigate what happens when the estimate of the gradient used in this section, contains additive noise.



Section 5.2.3      The Algorithm, Proof of Convergence, and Bounds on the Rate of Convergence if the Gradient is Estimated, and the Estimate is Noisy.

Using (5.1.8) the algorithm is

$$\underline{W}_{j+1} = \underline{W}_j - k \left\{ -2 (\underline{s}_j + \underline{n}_j) \left[ \underline{d}_j - (\underline{s}_j^T + \underline{n}_j^T) \underline{W}_j \right] + 2 K_1 \left[ \underline{W}_j^T \cdot \underline{n}_1 - a \right] \underline{n}_1 \right\} \quad (5.2.3.1)$$

These equations constitute a set of first-order stochastic difference equations. We will first solve for the asymptotic expected value of  $\underline{W}$ , which we will denote by  $\bar{\underline{W}}_\infty$ .

Taking the expected value of (5.2.3.1), under the assumption that  $E\{\underline{n}_j\} = \underline{0}$ ,  $E\{\underline{n}_j \underline{n}_j^T\} = \phi_n$ , and  $\underline{s}_j$ ,  $\underline{s}_k$ ,  $\underline{n}_l$ ,  $\underline{n}_m$  are statistically independent for  $k \neq j$  and  $n \neq m$ , we have

$$E\{\underline{W}_{j+1}\} = E\{\underline{W}_j\} - k \left\{ -2 \phi(\underline{s}, \underline{d}) + 2 \phi(\underline{s}, \underline{s}) E\{\underline{W}_j\} + 2 \phi_n E\{\underline{W}_j\} + 2 K_1 \left[ E\{\underline{W}_j^T\} \underline{n}_1 - a \right] \underline{n}_1 \right\}$$

Using (4.1.7a) yields

$$E\{\underline{W}_{j+1}\} = E\{\underline{W}_j\} + 2k \left[ \phi \underline{W}_{LMS} - \phi E\{\underline{W}_j\} - \phi_n E\{\underline{W}_j\} \right] - 2k K_1 \left[ E\{\underline{W}_j^T\} \underline{n}_1 - a \right] \underline{n}_1$$

We now have a set of deterministic first-order difference equations whose asymptotic value  $E\{\underline{W}\} \equiv \bar{\underline{W}}_\infty$ , can be found by setting  $E\{\underline{W}_j\} = E\{\underline{W}_{j+1}\} = \bar{\underline{W}}_\infty$ ; giving

$$\underline{W}_{LMS} - \bar{\underline{W}}_\infty - \phi^{-1} \phi_n \bar{\underline{W}}_\infty = K_1 \left[ \bar{\underline{W}}_\infty^T \cdot \underline{n}_1 - a \right] \phi^{-1} \underline{n}_1 \quad (5.2.3.2)$$

$$\text{let } \underline{\overline{W}}_{\infty} \equiv \underline{c} + d \underline{n}_1 \quad (5.2.3.3)$$

$$\text{where } \underline{c}^T \underline{n}_1 = \underline{n}_1^T \underline{c} = 0 \quad (5.2.3.4)$$

Remembering that  $\underline{n}_1^T \underline{n}_1 = 1$ , (5.2.3.2) becomes

$$\underline{W}_{LMS} - (I + \phi^{-1} \phi_n) (\underline{c} + d \underline{n}_1) = K_1 [d - a] \phi^{-1} \underline{n}_1 \quad (5.2.3.5)$$

Multiplying by  $\underline{n}_1^T (I + \phi^{-1} \phi_n)^{-1}$  on the left, and manipulating, gives

$$d = \frac{\underline{n}_1^T (I + \phi^{-1} \phi_n)^{-1} [\underline{W}_{LMS} + K_1 a \phi^{-1} \underline{n}_1]}{1 + K_1 \underline{n}_1^T (I + \phi^{-1} \phi_n)^{-1} \phi^{-1} \underline{n}_1} \quad (5.2.3.6)$$

From (5.2.3.5)

$$\underline{\overline{W}}_{\infty} = \underline{c} + d \underline{n}_1 = (I + \phi^{-1} \phi_n)^{-1} \left\{ \underline{W}_{LMS} - K_1 [d - a] \phi^{-1} \underline{n}_1 \right\}$$

Using (5.2.3.6), after some algebra, we get

$$\underline{\overline{W}}_{\infty} = \frac{(I + \phi^{-1} \phi_n)^{-1}}{1 + K_1 \underline{n}_1^T (I + \phi^{-1} \phi_n)^{-1} \phi^{-1} \underline{n}_1} \left\{ \left[ I + K_1 \underline{n}_1^T (I + \phi^{-1} \phi_n)^{-1} \phi^{-1} \underline{n}_1 I - K_1 \phi^{-1} \underline{n}_1 \underline{n}_1^T (I + \phi^{-1} \phi_n)^{-1} \right] \underline{W}_{LMS} + K_1 a \phi^{-1} \underline{n}_1 \right\} \quad (5.2.3.7)$$

If we let  $K_1 \rightarrow \infty$  we should get the same solution as equation (4.4.3.9), because the penalty function is infinite unless the weight vector lies exactly on the line  $\underline{W}^T \underline{n}_1 = a$ . Under these conditions, we get

$$\bar{W}_\infty = \frac{(I + \phi^{-1} \phi_n)^{-1}}{K_1 \underline{n}_1^T (I + \phi^{-1} \phi_n)^{-1} \phi^{-1} \underline{n}_1}$$

$$\left\{ \left[ K_1 \underline{n}_1^T (I + \phi^{-1} \phi_n)^{-1} \phi^{-1} \underline{n}_1 I - K_1 \phi^{-1} \underline{n}_1 \underline{n}_1^T (I + \phi^{-1} \phi_n)^{-1} \right] \underline{W}_{LMS} + K_1 a \phi^{-1} \underline{n}_1 \right\}$$

$$= (I + \phi^{-1} \phi_n)^{-1} \left\{ \underline{W}_{LMS} + \frac{-\phi^{-1} \underline{n}_1 \left\{ \underline{n}_1^T (I + \phi^{-1} \phi_n)^{-1} \underline{W}_{LMS} \right\} + a \phi^{-1} \underline{n}_1}{\underline{n}_1^T (I + \phi^{-1} \phi_n)^{-1} \phi^{-1} \underline{n}_1} \right\}$$

$$\therefore \bar{W}_\infty = (I + \phi^{-1} \phi_n)^{-1} \left\{ \underline{W}_{LMS} + \frac{\left[ a - \underline{n}_1^T (I + \phi^{-1} \phi_n)^{-1} \underline{W}_{LMS} \right] \phi^{-1} \underline{n}_1}{\underline{n}_1^T (I + \phi^{-1} \phi_n)^{-1} \phi^{-1} \underline{n}_1} \right\}$$

(5.2.3.8)

This is exactly the same as equation (4.3.3.9).

Again, because our difference equations describing the behavior of the weight vectors are stochastic, the above result is not sufficient to prove convergence of the weight vectors to  $\bar{W}_\infty$ , we must also show that the variance of the stochastic vectors  $\underline{q}_j = \underline{W}_j - \bar{W}_\infty$  is bounded. To do this

$$\text{define } \underline{q}_j \equiv \underline{W}_j - \bar{W}_\infty \quad (5.2.3.9)$$

the algorithm, (5.2.3.1) may be rewritten in the form

$$\begin{aligned} \underline{q}_{j+1} = \underline{q}_j - k \left[ 2(\underline{s}_j + \underline{n}_j) (\underline{s}_j^T + \underline{n}_j^T) + 2 K_1 \underline{n}_1 \underline{n}_1^T \right] \underline{q}_j \\ - k \left[ 2(\underline{s}_j + \underline{n}_j) (\underline{s}_j^T + \underline{n}_j^T) + 2 K_1 \underline{n}_1 \underline{n}_1^T \right] \bar{W}_\infty \\ + k \left[ 2(\underline{s}_j + \underline{n}_j) \underline{d}_j + 2 K_1 \underline{a} \underline{n}_1 \right] \end{aligned}$$

Define

$$\underline{H}_j \equiv 2(\underline{s}_j + \underline{n}_j) (\underline{s}_j^T + \underline{n}_j^T) + 2 K_1 \underline{n}_1 \underline{n}_1^T \quad (5.2.3.10)$$

$$\underline{V}_j \equiv 2(\underline{s}_j + \underline{n}_j) \underline{d}_j + 2 K_1 \underline{a} \underline{n}_1 \quad (5.2.3.11)$$

$$\therefore \underline{q}_{j+1} = \underline{q}_j - k \underline{\varphi}_j \quad (5.2.3.12)$$

where

$$\underline{\varphi}_j \equiv \underline{H}_j \underline{q}_j + \underline{h}_j \quad (5.2.3.13)$$

and

$$\underline{h}_j \equiv \underline{H}_j \bar{W}_\infty - \underline{V}_j \quad (5.2.3.14)$$

Note that  $E \{ H_j \}$  and  $E \{ \underline{h}_j \}$  are independent of  $j$ . Also  $H_j$  and  $\underline{h}_j$  are statistically independent of  $H_k$  and  $\underline{h}_k$  if  $j \neq k$  because we assumed  $\underline{s}_j, \underline{s}_k, \underline{n}_l, \underline{n}_m$  are statistically independent for  $k \neq j$  and  $n \neq m$ .

Again, as in the last section, it can be shown that

$$E \{ \underline{h}_j \} = \underline{0} \quad (5.2.3.15)$$

Note that  $E \{ H_j \} = 2(\phi + \phi_n) + 2K_1 \underline{n}_1 \underline{n}_1^T \equiv \underline{a}$  is a symmetric positive definite matrix.

The algorithm is thus

$$\underline{q}_{j+1} = \underline{q}_j - k \underline{\varphi}_j$$

where

$$\underline{\varphi}_j = H_j \underline{q}_j + \underline{h}_j$$

and  $H_j$  is a sequence of random  $n \times n$  matrices;  $\underline{h}_j$  is a sequence of random  $n$ -tuple vectors; the expected values of  $H_j$  and  $\underline{h}_j$  were shown to be independent of  $j$ ;  $H_j$  and  $\underline{h}_j$  are independent of  $H_l$  and  $\underline{h}_l$  for  $j \neq l$ ;  $E \{ \underline{h}_j \} = \underline{0}$ ; and the elements of  $H_j$  and  $\underline{h}_j$  have finite variance, with  $E \{ H_j \} = \underline{a}$ , where  $\underline{a}$  is a symmetric positive definite matrix.

Under these conditions, it is shown in Appendix A of chapter four, that for  $k$  sufficiently small

$$\lim_{j \rightarrow \infty} \| E \{ \underline{q}_j \} \| = 0 \quad (5.2.3.16)$$

and

$$\lim_{j \rightarrow \infty} \sup \| \underline{q}_j \| \leq V(k) \quad (5.2.3.17)$$

where the norm of a random vector  $\underline{u}$  is defined as

$$\| \underline{u} \| \equiv \sqrt{E \{ \underline{u}^T \underline{u} \}}$$

$$\text{and} \quad \lim_{k \rightarrow 0} V(k) = 0 \quad (5.2.3.18)$$

Equation (5.2.3.16) shows again that the random weight vectors converge, in the mean, to  $\bar{W}_\infty$  and (5.2.3.17) shows that the variance of the random weight vectors about their expected value is bounded, and the bound can be made as small as desired by choosing  $k$  sufficiently small as shown by (5.2.3.18).

The rate of convergence of the mean of the random weight vectors is shown in the proof of the above theorem to be bounded by  $\xi$ , where

$$\xi = || I - k (2\phi + 2\phi_n + 2K_1 \underline{n}_1 \underline{n}_1^T) || \quad (5.2.3.19)$$

Since  $A \equiv (2\phi + 2\phi_n + 2K_1 \underline{n}_1 \underline{n}_1^T)$  is positive definite symmetric, we have

$$\xi = \max \left\{ |1 - k\rho_1|, |1 - k\rho_n| \right\} \quad (5.2.3.20)$$

where  $\rho_1$  and  $\rho_n$  are both positive, and represent the minimum and maximum eigenvalues of  $A$  respectively, as shown by Goldstein<sup>(28)</sup> page 24.

Thus  $0 < \xi < 1$ .

In looking at the two approaches we have developed for adaptively optimizing the MSE subject to a constraint, the approach in chapter four represents an entirely new approach to the problem, whereas the approach in this chapter is essentially one of replacing the constrained problem by an unconstrained problem. Since stochastic unconstrained problems have already been well researched, we will not run computer simulations of the algorithm of this chapter, but will rather concentrate our efforts on the new algorithm developed in chapter four.

## CHAPTER 6

### Computer Simulations

In chapter three, we found the optimum SNR that we could achieve subject to a constraint on the super-gain ratio. Specifically, we showed that for a linear array of four isotropic detectors spaced  $d = 0.8\lambda$  ( $0.4\lambda$ ) apart, subject to the super-gain constraint  $Q = 0.08$  ( $0.11$ ), embedded in a uniform noise field, with a normalized signal impinging from broadside (endfire), the best SNR we could get at the array output was  $0.187$  ( $0.438$ ).

In this chapter we will simulate a projected gradient algorithm which automatically makes an array of four isotropic detectors spaced  $d = 0.8\lambda$  ( $0.4\lambda$ ) apart maximize the average output SNR, subject to the constraint that the super-gain ratio  $Q$  is  $\leq 0.08$  ( $0.11$ ) when the signal impinges from broadside (endfire) and the noise is isotropic.

We will again (as in chapter three) assume that the signal and noise are sufficiently temporally narrowband so that the filter following each detector can be implemented by only two taps (or attenuators) separated by a quarter period delay as shown in Fig. 6.2.1 when using the multichannel filter point of view. This corresponds to Fig. 3.1.2 when using the antenna point of view.

We will formulate the problem first from the antenna point of view, i.e. we will write the SNR and super-gain ratio (Q-factor) in terms of the real and imaginary parts of the detector currents  $I_{1r}$ ,  $I_{1i}$ ,  $I_{2r}$ ,  $I_{2i}$ , ...,  $I_{4r}$ ,  $I_{4i}$ , and second from the multichannel filter point of view, i.e. we will write the SNR in terms of  $w_1$ ,  $w_2$ , ...,  $w_7$ ,  $w_8$ . In agreement with the results of chapter two, we will observe that  $I_{1r}$  is equivalent to  $w_1$ ,  $I_{1i}$  is equivalent to  $w_2$ ,  $I_{2r}$  is equivalent to  $w_3$ , etc. We will then use this equivalence to write the expression for the super-gain ratio in terms of  $w_1, \dots, w_8$ .

## Section 6.1 Antenna Theory Approach

When the signal is impinging from broadside, the time average signal power coming out of the array is given by equation (2.1.10)

$$S = \frac{1}{2} \underline{I}^* \underline{V}_1 \underline{V}_1^* \underline{I}$$

where  $\underline{V}_1$  is given by equation (3.1.14)

$$\underline{V}_1 = \text{col} [1 \ 1 \ 1 \ 1]$$

Writing  $\underline{I}$  as  $\text{col} [I_{1r} + jI_{1i}, I_{2r} + jI_{2i}, I_{3r} + jI_{3i}, I_{4r} + jI_{4i}]$ ,

expanding and then rearranging gives

$$S = \frac{1}{2} \underline{I}^T \begin{bmatrix} 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \end{bmatrix} \underline{I} \quad (6.1.1)$$

where now  $\underline{I} = \text{col} [I_{1r}, I_{1i}, I_{2r}, I_{2i}, I_{3r}, I_{3i}, I_{4r}, I_{4i}]$  (6.1.2)

Assuming the noise field is uniform, as in chapter three, the time average noise power coming out of the array is given by equation (2.1.12)

$$N = \frac{1}{2} \underline{I}^* A \underline{I}$$

where  $A$  is given by (3.1.11).

This expression can be manipulated into

$$N = \frac{1}{2} \underline{I}^T E \underline{I} \quad (6.1.3)$$

where  $\underline{I}$  is given by (6.1.2) and



E=

$$\begin{bmatrix}
 4\pi & 0 & \frac{2\lambda}{d} \sin \frac{2\pi d}{\lambda} & 0 & \frac{\lambda}{d} \sin \frac{4\pi d}{\lambda} & 0 & \frac{2\lambda}{3d} \sin \frac{6\pi d}{\lambda} & 0 \\
 0 & 4\pi & 0 & \frac{2\lambda}{d} \sin \frac{2\pi d}{\lambda} & 0 & \frac{\lambda}{d} \sin \frac{4\pi d}{\lambda} & 0 & \frac{2\lambda}{3d} \sin \frac{6\pi d}{\lambda} \\
 \frac{2\lambda}{d} \sin \frac{2\pi d}{\lambda} & 0 & 4\pi & 0 & \frac{2\lambda}{d} \sin \frac{2\pi d}{\lambda} & 0 & \frac{\lambda}{d} \sin \frac{4\pi d}{\lambda} & 0 \\
 0 & \frac{2\lambda}{d} \sin \frac{2\pi d}{\lambda} & 0 & 4\pi & 0 & \frac{2\lambda}{d} \sin \frac{2\pi d}{\lambda} & 0 & \frac{\lambda}{d} \sin \frac{4\pi d}{\lambda} \\
 \frac{\lambda}{d} \sin \frac{4\pi d}{\lambda} & 0 & \frac{2\lambda}{d} \sin \frac{2\pi d}{\lambda} & 0 & 4\pi & 0 & \frac{2\lambda}{d} \sin \frac{2\pi d}{\lambda} & 0 \\
 0 & \frac{\lambda}{d} \sin \frac{4\pi d}{\lambda} & 0 & \frac{2\lambda}{d} \sin \frac{2\pi d}{\lambda} & 0 & 4\pi & 0 & \frac{2\lambda}{d} \sin \frac{2\pi d}{\lambda} \\
 \frac{2\lambda}{3d} \sin \frac{6\pi d}{\lambda} & 0 & \frac{\lambda}{d} \sin \frac{4\pi d}{\lambda} & 0 & \frac{2\lambda}{d} \sin \frac{2\pi d}{\lambda} & 0 & 4\pi & 0 \\
 0 & \frac{2\lambda}{3d} \sin \frac{6\pi d}{\lambda} & 0 & \frac{\lambda}{d} \sin \frac{4\pi d}{\lambda} & 0 & \frac{2\lambda}{d} \sin \frac{2\pi d}{\lambda} & 0 & 4\pi
 \end{bmatrix}$$

(6.1.4)

In terms of this eight dimensional  $\underline{I}$  vector, the Q factor is given by (see equation (3. 1. 13))

$$Q = \frac{\underline{I}^T \underline{I}}{\underline{I}^T \underline{E} \underline{I}} \quad (6. 1. 5)$$

If the signal impinges from endfire the only quantity that changes in the above formulation is the time average signal power S (the noise power and the Q factor are the same as for the broadside signal case). Now

$$S = \frac{1}{2} \underline{I}^* \underline{V}_1 \underline{V}_1^* \underline{I}$$

where from (3. 1. 15),

$$\underline{V}_1 = \text{col} \begin{bmatrix} e^{j(-\frac{3\pi d}{\lambda})} & e^{j(-\pi \frac{d}{\lambda})} & e^{j(\pi \frac{d}{\lambda})} & e^{j(3\pi \frac{d}{\lambda})} \end{bmatrix}$$

This expression may be manipulated into

$$S = \frac{1}{2} \underline{I}^T \underline{F} \underline{I} \quad (6. 1. 6)$$

where  $\underline{I}$  is given by (6. 1. 2) and

$$F = \begin{bmatrix} 1 & 0 & \cos \frac{2\pi d}{\lambda} & \sin \frac{2\pi d}{\lambda} & \cos \frac{4\pi d}{\lambda} & \sin \frac{4\pi d}{\lambda} & \cos \frac{6\pi d}{\lambda} & \sin \frac{6\pi d}{\lambda} \\ 0 & 1 & -\sin \frac{2\pi d}{\lambda} & \cos \frac{2\pi d}{\lambda} & -\sin \frac{4\pi d}{\lambda} & \cos \frac{4\pi d}{\lambda} & -\sin \frac{6\pi d}{\lambda} & \cos \frac{6\pi d}{\lambda} \\ \cos \frac{2\pi d}{\lambda} & -\sin \frac{2\pi d}{\lambda} & 1 & 0 & \cos \frac{2\pi d}{\lambda} & \sin \frac{2\pi d}{\lambda} & \cos \frac{4\pi d}{\lambda} & \sin \frac{4\pi d}{\lambda} \\ \sin \frac{2\pi d}{\lambda} & \cos \frac{2\pi d}{\lambda} & 0 & 1 & -\sin \frac{2\pi d}{\lambda} & \cos \frac{2\pi d}{\lambda} & -\sin \frac{4\pi d}{\lambda} & \cos \frac{4\pi d}{\lambda} \\ \cos \frac{4\pi d}{\lambda} & -\sin \frac{4\pi d}{\lambda} & \cos \frac{2\pi d}{\lambda} & -\sin \frac{2\pi d}{\lambda} & 1 & 0 & \cos \frac{2\pi d}{\lambda} & \sin \frac{2\pi d}{\lambda} \\ \sin \frac{4\pi d}{\lambda} & \cos \frac{4\pi d}{\lambda} & \sin \frac{2\pi d}{\lambda} & \cos \frac{2\pi d}{\lambda} & 0 & 1 & -\sin \frac{2\pi d}{\lambda} & \cos \frac{2\pi d}{\lambda} \\ \cos \frac{6\pi d}{\lambda} & -\sin \frac{6\pi d}{\lambda} & \cos \frac{4\pi d}{\lambda} & -\sin \frac{4\pi d}{\lambda} & \cos \frac{2\pi d}{\lambda} & -\sin \frac{2\pi d}{\lambda} & 1 & 0 \\ \sin \frac{6\pi d}{\lambda} & \cos \frac{6\pi d}{\lambda} & \sin \frac{4\pi d}{\lambda} & \cos \frac{4\pi d}{\lambda} & \sin \frac{2\pi d}{\lambda} & \cos \frac{2\pi d}{\lambda} & 0 & 1 \end{bmatrix}$$

(6.1.7).

## Section 6.2 Multichannel Filter Approach

Let us now find the time average output signal power due to a deterministic signal generated by a far field point source (see Fig. 6.2.1).

At each detector, the signal is given by

$$\begin{aligned} & \text{Re } e^{-j \frac{2\pi}{\lambda} \underline{u}_0 \cdot \underline{r}_i} e^{j\omega t} \\ &= \cos \left( \omega t - \frac{2\pi}{\lambda} \underline{u}_0 \cdot \underline{r}_i \right) \end{aligned} \quad (6.2.1)$$

The output  $y(t)$  due to the signal is

$$y(t) = [w_1 \ w_2 \ w_3 \ w_4 \ w_5 \ w_6 \ w_7 \ w_8] \begin{bmatrix} \cos \left( \omega t - \frac{2\pi}{\lambda} \underline{u}_0 \cdot \underline{r}_1 \right) \\ \cos \left( \omega t - \frac{2\pi}{\lambda} \underline{u}_0 \cdot \underline{r}_1 - \omega \Delta \right) \\ \cos \left( \omega t - \frac{2\pi}{\lambda} \underline{u}_0 \cdot \underline{r}_2 \right) \\ \cos \left( \omega t - \frac{2\pi}{\lambda} \underline{u}_0 \cdot \underline{r}_2 - \omega \Delta \right) \\ \cos \left( \omega t - \frac{2\pi}{\lambda} \underline{u}_0 \cdot \underline{r}_3 \right) \\ \cos \left( \omega t - \frac{2\pi}{\lambda} \underline{u}_0 \cdot \underline{r}_3 - \omega \Delta \right) \\ \cos \left( \omega t - \frac{2\pi}{\lambda} \underline{u}_0 \cdot \underline{r}_4 \right) \\ \cos \left( \omega t - \frac{2\pi}{\lambda} \underline{u}_0 \cdot \underline{r}_4 - \omega \Delta \right) \end{bmatrix} \equiv \underline{W}^T \underline{a}_1 \quad (6.2.2)$$

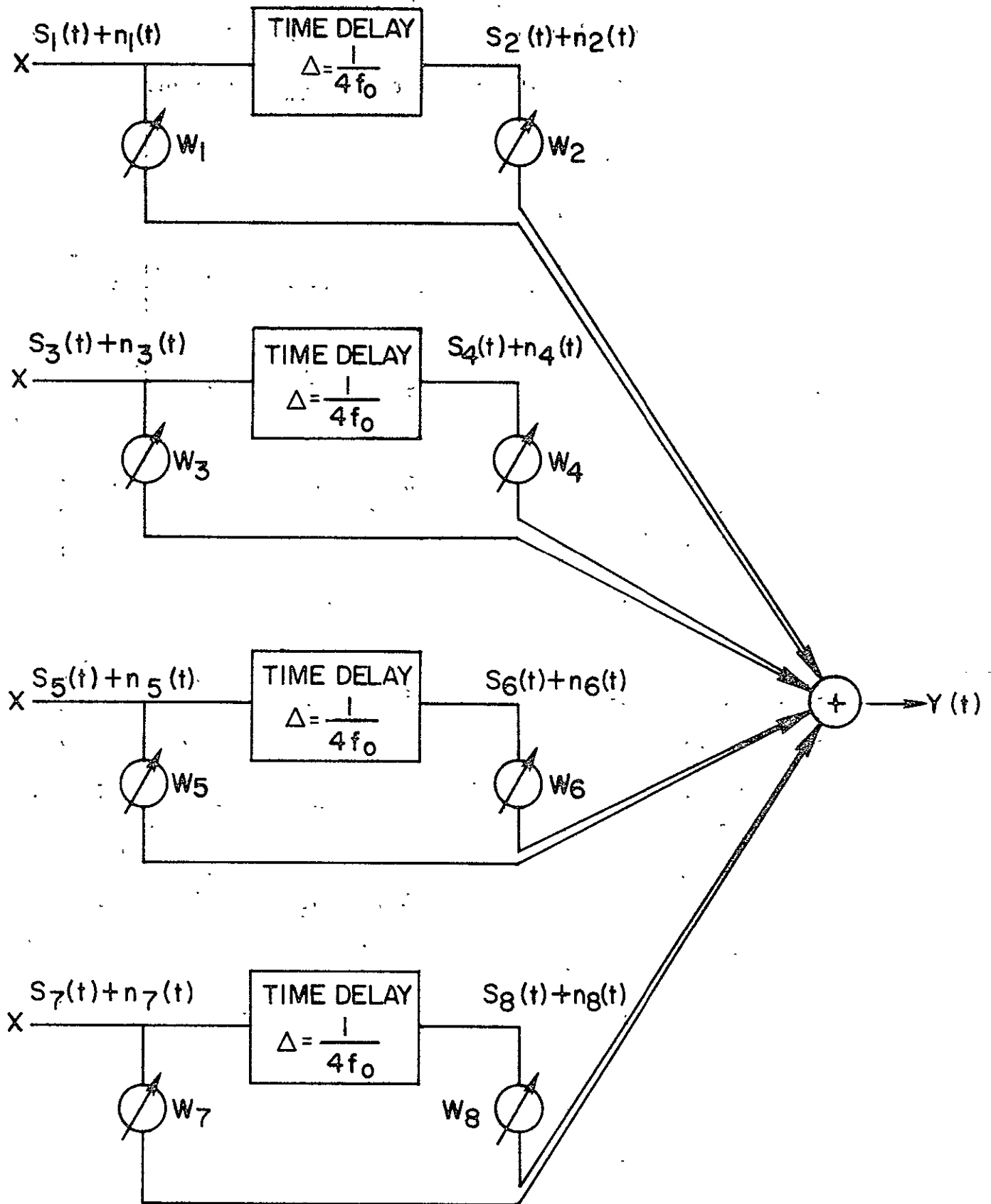


Fig. 6.2.1 Processor Structure

The signal output power is  $S(t) = \underline{W}^T \underline{a}_1 \underline{a}_1^T \underline{W}$  (6.2.3)

For the case of a broadside signal,  $\underline{u}_0 \cdot \underline{r}_1 = \underline{u}_0 \cdot \underline{r}_2 = \underline{u}_0 \cdot \underline{r}_3 = \underline{u}_0 \cdot \underline{r}_4 = 0$ .

Letting  $d = \cos \omega t$ ,  $e = \cos (\omega t - \omega \Delta)$ , the matrix  $\underline{a}_1 \underline{a}_1^T$  is given by

$$\underline{a}_1 \underline{a}_1^T = \begin{bmatrix} d^2 & de & d^2 & de & d^2 & de & d^2 & de \\ ed & e^2 & ed & e^2 & ed & e^2 & ed & e^2 \\ d^2 & de & d^2 & de & d^2 & de & d^2 & de \\ ed & e^2 & ed & e^2 & ed & e^2 & ed & e^2 \\ d^2 & de & d^2 & de & d^2 & de & d^2 & de \\ ed & e^2 & ed & e^2 & ed & e^2 & ed & e^2 \\ d^2 & de & d^2 & de & d^2 & de & d^2 & de \\ ed & e^2 & ed & e^2 & ed & e^2 & ed & e^2 \end{bmatrix}$$

Since  $\frac{1}{2\pi} \int_0^{2\pi/\omega} \cos^2 \omega t dt = \frac{1}{2}$  and  $\frac{1}{2\pi} \int_0^{2\pi/\omega} \cos \omega t \cos (\omega t - \omega \Delta) dt = \frac{1}{2} \cos \omega \Delta$

the time average signal power output is given by  $\underline{W}^T R \underline{W}$  where the matrix R is

$$R = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta \\ \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta \\ \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta \\ \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta \\ \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} & \frac{1}{2} \cos \omega \Delta & \frac{1}{2} \end{bmatrix}$$

Since  $\cos \omega \Delta = \cos \frac{\pi}{2} = 0$ , this expression for the average signal power becomes identical to equation (6.1.1) with the vector  $\underline{W}$  replacing the vector  $\underline{I}$  of equation (6.1.2). Similarly we can show that the expressions representing the time average noise power in terms of  $\underline{I}$  and  $\underline{W}$  are identical if we replace  $\underline{I}$  by  $\underline{W}$ , i.e.

$$N = \frac{1}{2} \underline{W}^T E \underline{W} \quad (6.2.5)$$

where  $E$  is given by (6.1.4)

$$Q = \frac{\underline{W}^T \underline{W}}{\underline{W}^T E \underline{W}} \quad (6.2.6)$$

and, if the signal impinges from endfire

$$S = \frac{1}{2} \underline{W}^T F \underline{W} \quad (6.2.7)$$

where  $F$  is given by (6.1.7)



### Section 6.3 Maximization of SNR Subject to $Q \leq q$

The reason we went through two separate formulations of the same physical problem in sections 6.1 and 6.2 is as follows: In the  $\underline{W}$  formulation the numerator matrix in the expression for the SNR is of rank two, and this makes it impossible for us to conclude from this formulation that the SNR is a concave function of the  $\underline{W}$ 's, and hence possesses a unique maximum. However by using the complex  $\underline{I}$  formulation, we will be able to show that there exists one unique value of  $\underline{I}$  (and hence by our analogy, one unique value of  $\underline{W}$ ) which maximizes the SNR. The proof is as follows:

By equation (2.1.13)

$$\text{SNR} = \frac{\underline{I}^* \underline{V}_1 \underline{V}_1^* \underline{I}}{\underline{I}^* \underline{A} \underline{I}} \quad (6.3.1)$$

Let us take the first variation of the SNR with respect to the complex vector  $\underline{I}$  and set it equal to zero to find the possible extreme points.

$$\begin{aligned} \delta(\text{SNR}) &= \frac{(\underline{I}^* \underline{A} \underline{I}) \left[ \underline{I}^* \underline{V}_1 \underline{V}_1^* (\delta \underline{I}) + (\delta \underline{I}^*) \underline{V}_1 \underline{V}_1^* \underline{I} \right]}{(\underline{I}^* \underline{A} \underline{I})^2} \\ &\quad - \frac{(\underline{I}^* \underline{V}_1 \underline{V}_1^* \underline{I}) \left[ \underline{I}^* \underline{A} (\delta \underline{I}) + (\delta \underline{I}^*) \underline{A} \underline{I} \right]}{(\underline{I}^* \underline{A} \underline{I})^2} = 0 \end{aligned} \quad (6.3.2)$$

$$\text{Letting } \underline{y}^* \equiv (\underline{I}^* \underline{A} \underline{I}) \underline{I}^* \underline{V}_1 \underline{V}_1^* - (\underline{I}^* \underline{V}_1 \underline{V}_1^* \underline{I}) \underline{I}^* \underline{A} \quad (6.3.3)$$

equation (6.3.2) becomes after rearranging

$$\underline{y}^* \delta \underline{I} + \delta \underline{I}^* \underline{y} = 0 \quad (6.3.4)$$

Since this equation must hold for arbitrary  $\delta \underline{I}$  where  $\underline{I}$  is complex, (6.3.4) implies that  $\underline{y} = \underline{0}$ , which implies

$$(\underline{I}^* \underline{A} \underline{I}) \underline{V}_1 (\underline{V}_1^* \underline{I}) = (\underline{I}^* \underline{V}_1 \underline{V}_1^* \underline{I}) \underline{A} \underline{I}$$

$$(\underline{I}^* \underline{A} \underline{I}) \underline{V}_1 = (\underline{I}^* \underline{V}_1) \underline{A} \underline{I} \quad (6.3.5)$$

This equation is satisfied if  $(\underline{I}^* \underline{V}_1) = 0$ , which would mean that equation (6.3.1) was zero, obviously a minimum value, or if

$$\underline{I} = \frac{(\underline{I}^* \underline{A} \underline{I})}{(\underline{I}^* \underline{V}_1)} \underline{A}^{-1} \underline{V}_1 \quad (6.3.6)$$

This value of  $\underline{I}$  gives the unique maximum of the SNR.

There is also only one unique minimum. Corresponding to these two values of  $\underline{I}$ , there is a unique value of  $\underline{W}$  which maximizes the SNR, and one unique value which minimizes the SNR.

It is easy to prove that the set of points  $\underline{W}$  which satisfy  $Q(\underline{W}) \leq q$  is star connected about  $\underline{W}_0 = \underline{0}$ , by observing that if  $Q(\underline{W}) \leq q$ , then  $Q(\underline{x})$  where  $\underline{x} = \lambda \underline{W} + (1 - \lambda) \underline{W}_0$ ,  $0 \leq \lambda \leq 1 \implies \underline{x} = \lambda \underline{W}$  also satisfies  $Q(\underline{x}) \leq q$ . This star connectedness is a consequence of the fact that the Q factor is independent of the magnitude of  $\underline{W}$ .

Because the region  $Q(\underline{W}) \leq q$  is connected and the objective function  $SNR(\underline{W})$  is concave, our projection algorithm will converge to the constrained maximum, which occurs at the unconstrained maximum of the SNR, or on the boundary of the feasible region (in the broadside and endfire cases under study, we know that the unconstrained maximum of the SNR lies outside the feasible region by the graphs in chapter three).

Since the solution to the problem of maximizing the SNR subject to the constraint  $Q \leq q$ , lies on the boundary (i.e.  $Q = q$ ), the Lagrange solution we found in chapter three is also the solution we should wind up with in this chapter.

## Section 6.4 The Gradient Projection Algorithm

The function to be maximized is  $SNR = \frac{\underline{W}^T \underline{F} \underline{W}}{\underline{W}^T \underline{E} \underline{W}}$  subject to the constraint  $Q = \frac{\underline{W}^T \underline{W}}{\underline{W}^T \underline{A} \underline{W}} \leq q_0$ . Note that since the signal direction is assumed known to us (i.e.  $\underline{F}$  is known), we never need to know the signal itself (as opposed to needing  $\underline{d}_j$  when we used a MSE criterion in chapter four). We will investigate three cases:

1. The spatial distribution of the noise is known a priori (i.e. the elements of the matrix  $\underline{E}$  are known) and there is no additive self-noise associated with each detector.
2. The spatial distribution of the noise is unknown (i.e.  $\underline{E}$  must be estimated from observations of the detector outputs when there is no signal present) and there is no additive self-noise associated with each detector.
3. The spatial distribution of the noise is unknown and there is additive self-noise associated with each detector.

Before we describe the algorithm, note that the gradient of the SNR is given by

$$\nabla_{\underline{W}} (SNR) = \frac{-(\underline{W}^T \underline{F} \underline{W})^2 \underline{E} \underline{W} + (\underline{W}^T \underline{E} \underline{W})^2 \underline{F} \underline{W}}{(\underline{W}^T \underline{E} \underline{W})^2} \quad (6.4.1)$$

Also note that the normal to the hyperplane tangent to the surface

$$Q = \frac{\underline{W}^T \underline{W}}{\underline{W}^T \underline{A} \underline{W}} = q_0 \text{ is given by}$$

$$\underline{n} = \frac{-(\underline{W}^T \underline{W})^2 \underline{A} \underline{W} + (\underline{W}^T \underline{A} \underline{W})^2 \underline{W}}{(\underline{W}^T \underline{A} \underline{W})^2} \quad (6.4.2)$$

Our algorithm works as follows: We start at any arbitrary value  $\underline{W}_0$  ( $\underline{W}_0 = \text{col} [11111111]$ ). We check to see if  $\underline{W}_0$  satisfies the constraint, (if it does not, we keep moving in the direction  $-\underline{n}$ , i.e.  $\underline{W}_{i+1} = \underline{W}_i - k\underline{n}$ , until we arrive at a value of  $\underline{W}$  which does satisfy the constraint). In case 1, we try to move in the direction given by the gradient, i.e.

$$\underline{W}_{j+1} = \underline{W}_j + k \nabla_{\underline{W}}(\text{SNR}) \quad (6.4.3)$$

where  $\nabla_{\underline{W}}(\text{SNR})$  is given by (6.4.1). We next check  $\underline{W}_{j+1}$  to make sure it satisfies the constraint. If it does, we continue our iterations as given by equation (6.4.3) indefinitely. If, on the other hand  $\underline{W}_{j+1}$  does not satisfy the constraint, we form a different  $\underline{W}_{j+1}$  given by

$$\underline{W}_{j+1} = \underline{W}_j + k P \nabla_{\underline{W}}(\text{SNR}) \quad (6.4.4)$$

where  $P$ , the projection matrix, is given by  $I - \underline{n} \underline{n}^T$  and  $\underline{n}$  is given by (6.4.2). Provided  $k$  is "small enough," this value of  $\underline{W}_{j+1}$  will always satisfy the constraint and give a higher value of SNR than  $\underline{W}_j$ , because we are projecting the gradient into the hyperplane tangent to the constraint as shown in Fig. 6.4.1.

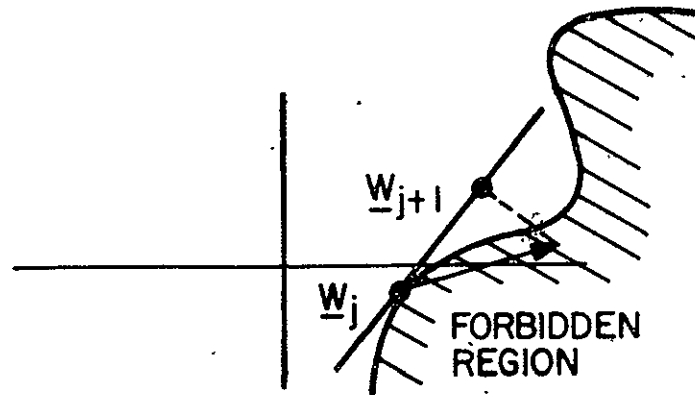


Fig. 6.4.1 Gradient Projection Operation

The reason  $k$  must be "small enough" is intuitively clear from the same figure. If we move too far along the hyperplane tangent to the constraint at  $\underline{w}_j$ , we may not satisfy the constraint at  $\underline{w}_{j+1}$ . In order to resolve this problem in our simulations, we chose  $k$  so as to make the square of the norm of  $k \nabla_{\underline{w}} (\text{SNR})$  equal to 0.001 times the square of the norm of  $\underline{w}_j$ , i.e.

$$k = \sqrt{0.001} \frac{\|\underline{w}_j\|}{\|\nabla_{\underline{w}}(\text{SNR})\|} \quad (6.4.5)$$

In case 2 where the noise correlation matrix  $E$  is unknown, for each element  $E_{ij} = E \{n_i(t) n_j(t)\}$  of the matrix  $E$  we substituted the instantaneous value of the correlation, i.e.  $E \rightarrow \tilde{E}$  where  $\tilde{E}_{ij}$  at iteration  $k$  is given by  $n_i(t_k) n_j(t_k)$  (see Fig. 6.2.1). In chapter four we proved that we would get convergence by using this substitution if our criterion was to minimize the MSE subject to a linear constraint.

In case 3 we substituted the matrix  $\tilde{E}$  for the matrix  $E$  in (6.4.1) where  $\tilde{E}_{ij}$  at iteration  $k$  is given by  $\tilde{E}_{ij} = [n_i(t_k) + \xi_i(t_k)] [n_j(t_k) + \xi_j(t_k)]$  where  $\xi(t_k)$  is white gaussian noise of variance 0.1.

To generate the vector random variables  $\underline{n}_k$  such that  $E\{\underline{n}_k \underline{n}_k^T\} \equiv E$ , we did the following:  $E$  is a positive definite matrix so that it possesses a square root, call the square root matrix  $E^{\frac{1}{2}}$ , where  $E^{\frac{1}{2}} E^{\frac{1}{2}} = E$ . We generated a vector random variable  $\underline{V}$ , all of whose components were zero mean independent gaussian random variables with variance one. Then

$$\underline{N}_k = E^{\frac{1}{2}} \underline{V} \quad (6.4.6)$$

and  $\underline{n}_k$  satisfies  $E\{\underline{n}_k \underline{n}_k^T\} = E\{E^{\frac{1}{2}} \underline{V} \underline{V}^T E^{\frac{1}{2}}\}$

$$= E^{\frac{1}{2}} E\{\underline{V} \underline{V}^T\} E^{\frac{1}{2}} = E^{\frac{1}{2}} I E^{\frac{1}{2}} = E \text{ as required.}$$

We simulated the aforementioned three cases for a signal impinging from both broadside and endfire and obtained the results shown in Figures 6.4.2 - 6.4.7. Note that in case 1 where the E matrix, and hence the gradient, was known we used  $k = 0.5$ , and we did not normalize  $k$  by equation (6.4.5).

By comparing Figs. 6.4.3 to 6.4.4 and 6.4.6 to 6.4.7, it can be seen that, as expected, the algorithm converges to the constrained optimal value faster, and there is less variance about the optimal value, when there is no additive detector noise present.

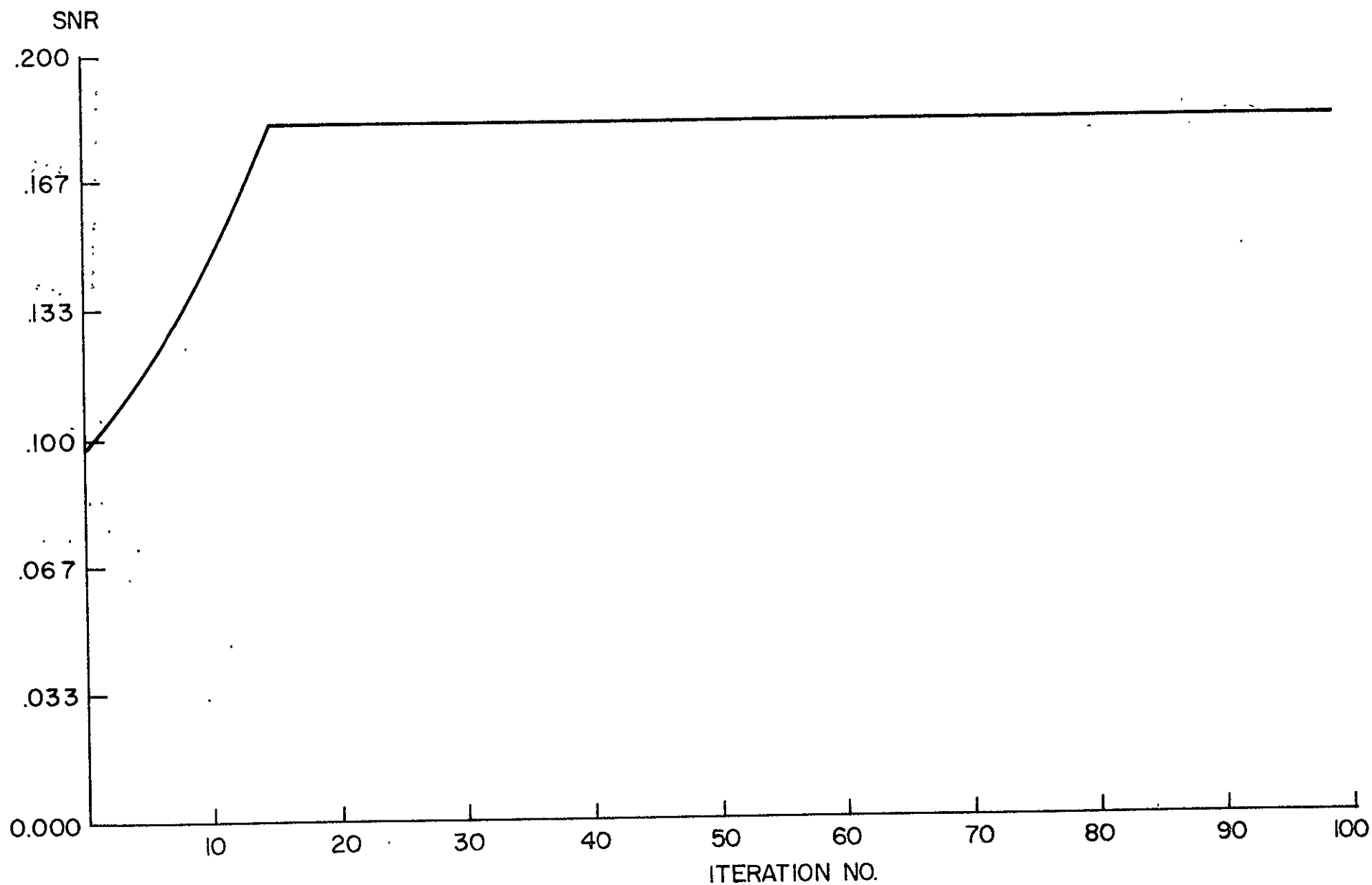


Fig. 6.4.2. Broadside. Gradient Known, No Additive Detector Noise

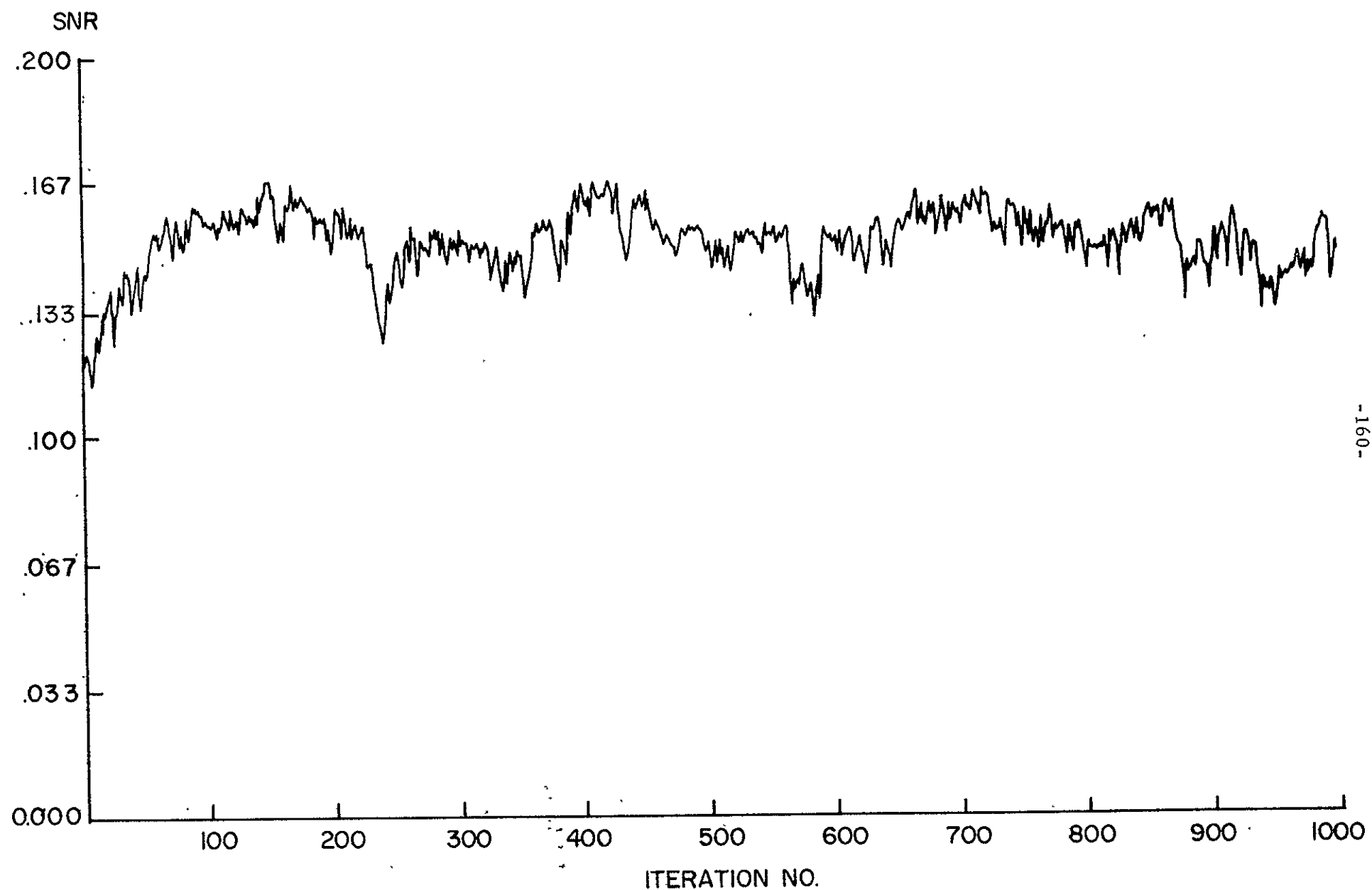


Fig. 6.4.3. Broadside. Gradient Estimated, No Additive Detector Noise



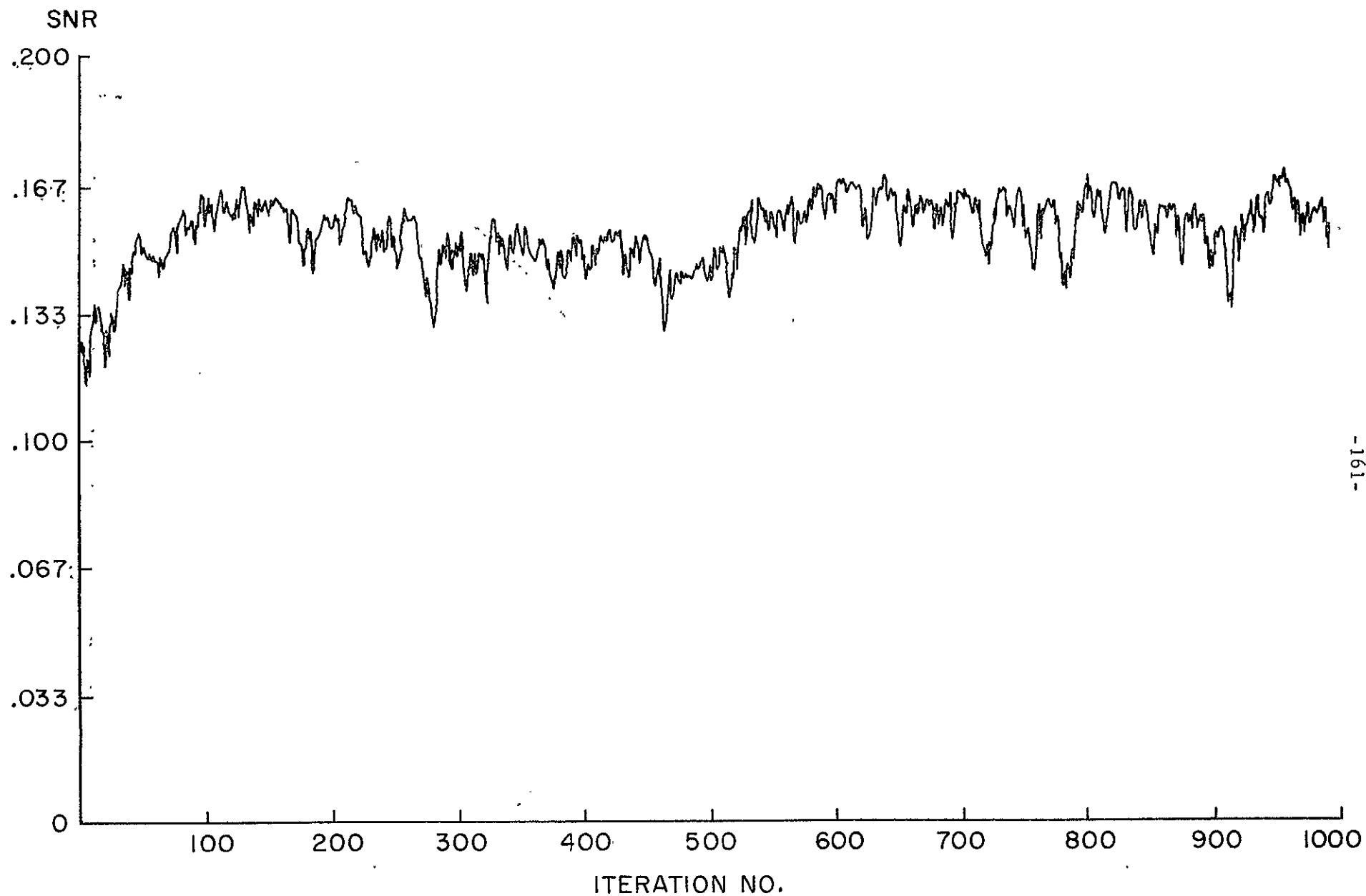


Fig. 6.4.4. Broadside. Gradient Estimated, Plus Additive Detector Noise

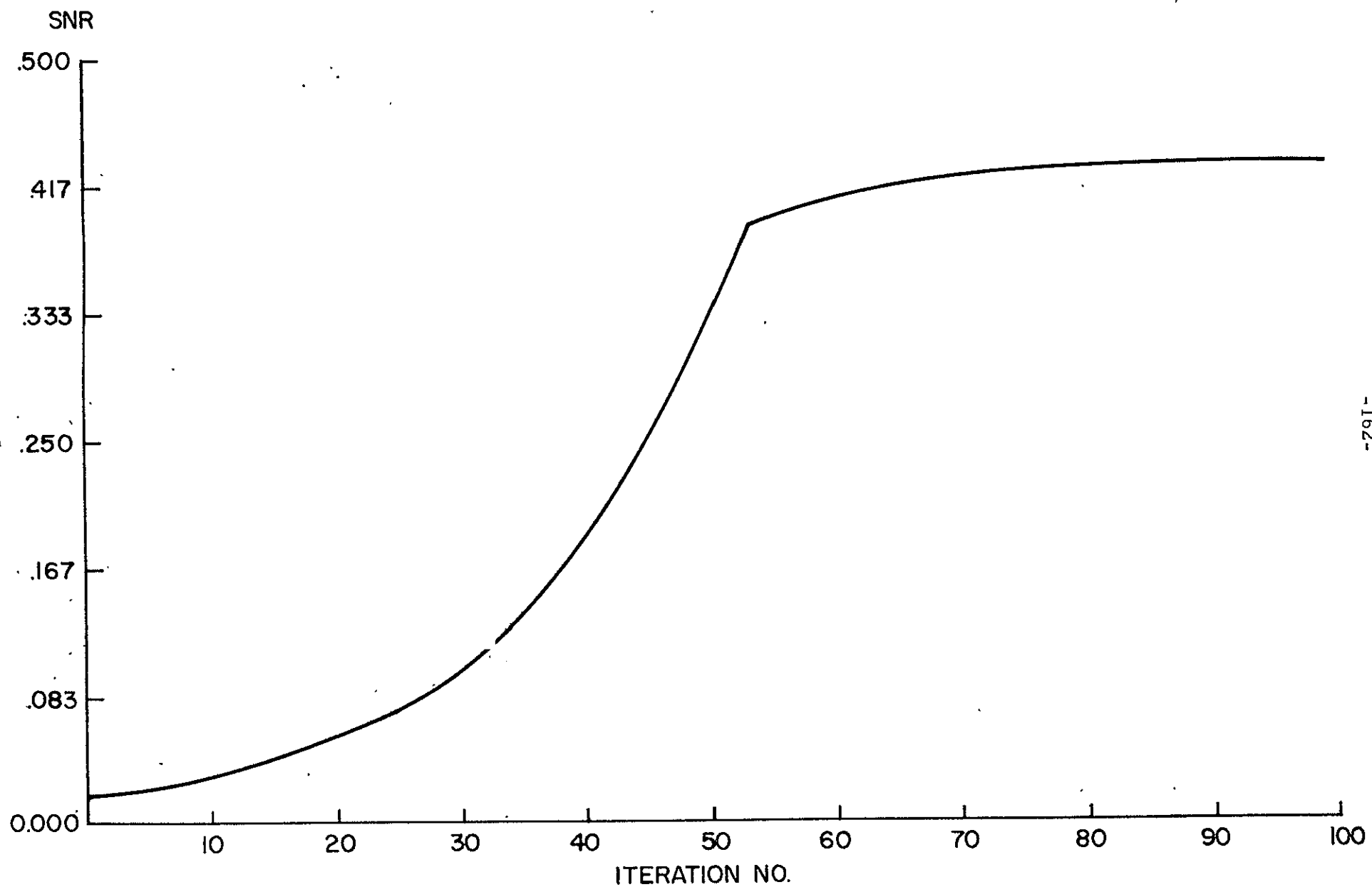


Fig. 6.4.5. Endfire. Gradient Known, No Additive Detector Noise

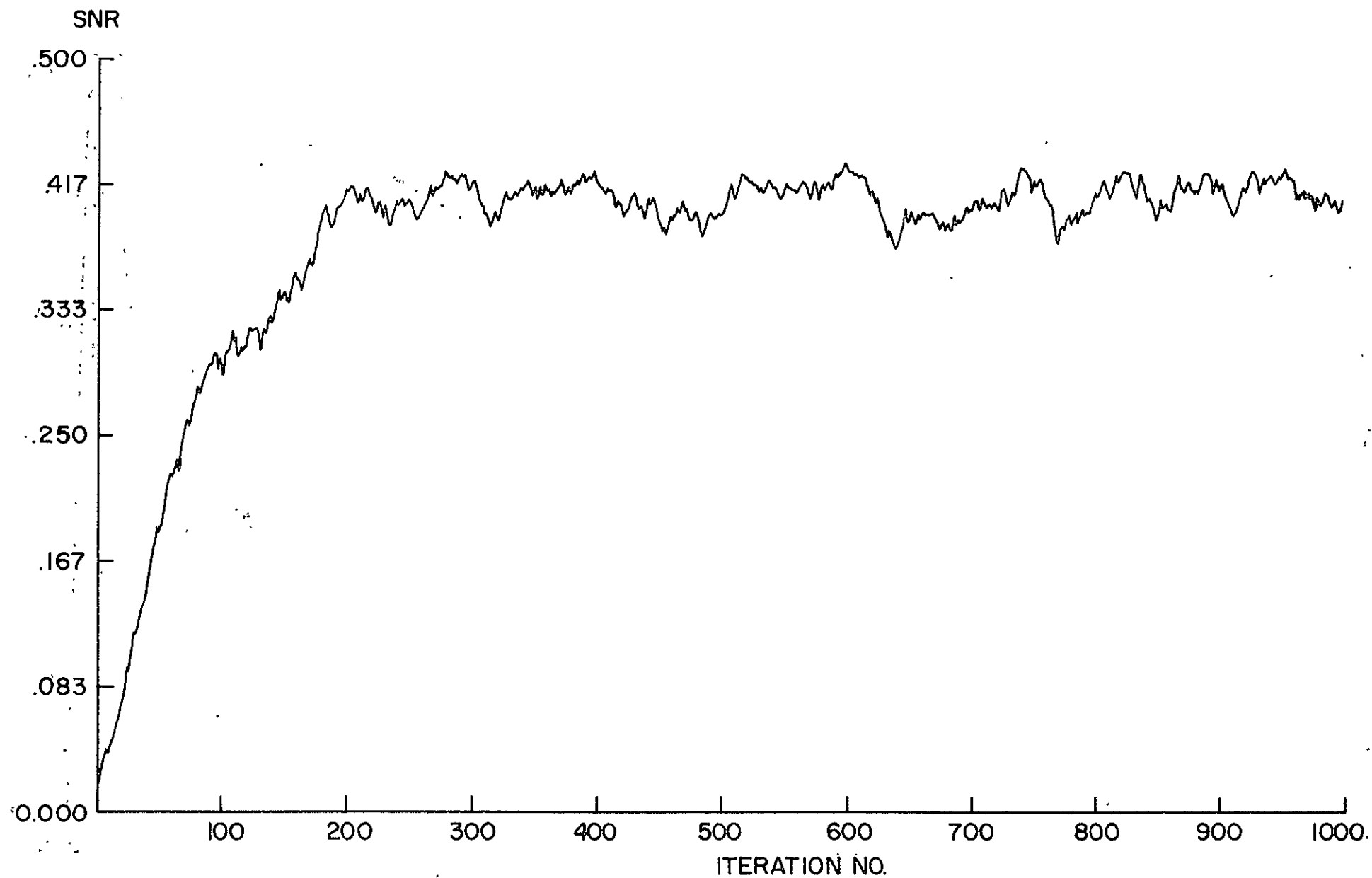


Fig. 6.4.6. Endfire. Gradient Estimated, No Additive Detector Noise

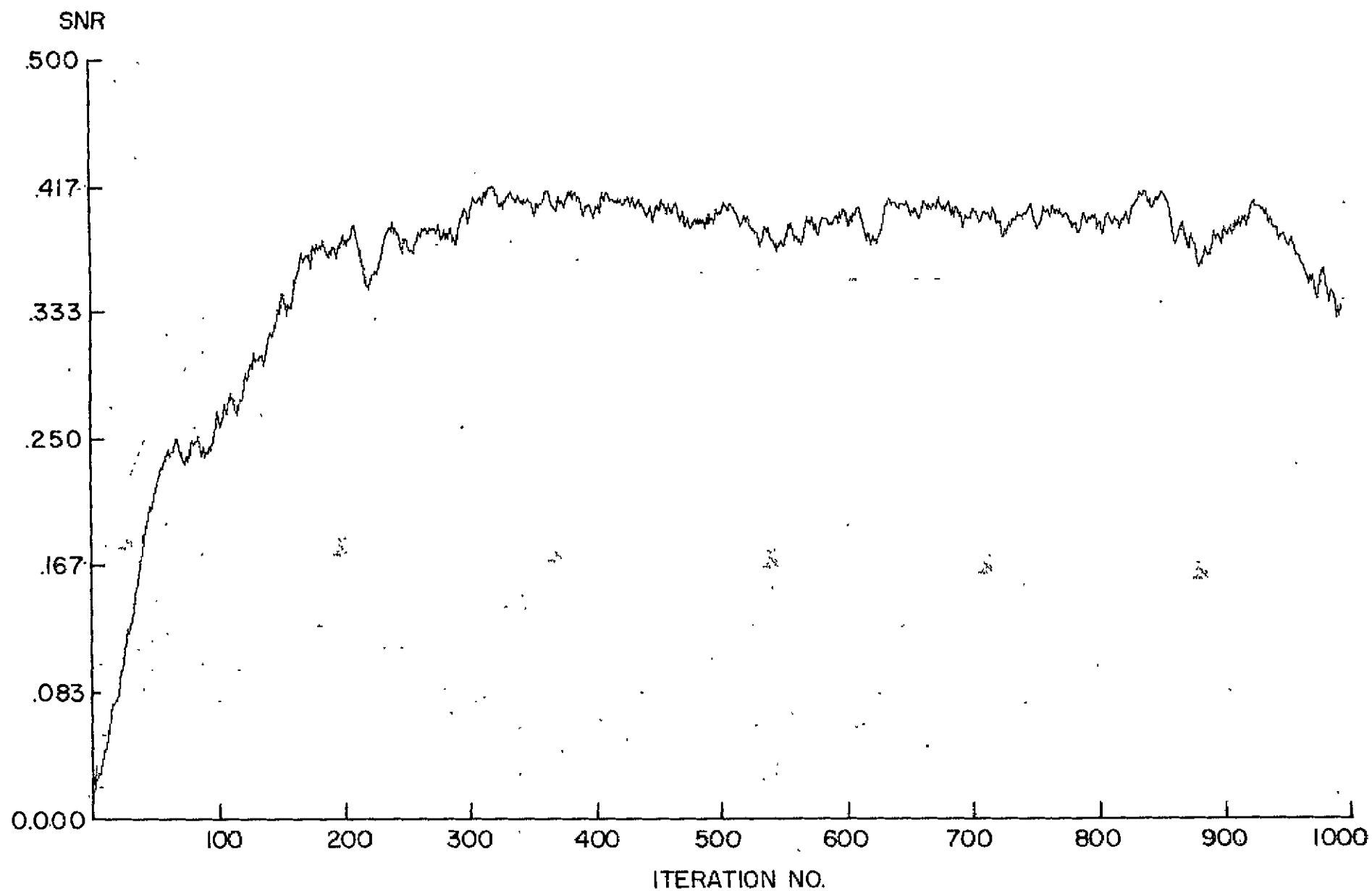


Fig. 6.4.7. Endfire. Gradient Estimated, Plus Additive Detector Noise

## Section 6.5 Conclusions

We have presented and analyzed two stochastic gradient algorithms, which can be used to find a constrained optimum point for a concave or convex objective function subject to constraints which form a connected region, even when we do not have the objective function available, but only have a noisy estimate of the objective function. When the constraints consisted of only one linear constraint, we proved convergence to the constrained optimum value and bounded the rate of convergence of the algorithms to the constrained optimum value.

## REFERENCES

1. Taçk, D. H., "Correlation Detection Arrays Designed for Directional Noise Fields," J. Acoustical Soc. Am., vol. 42, No. 3, pp. 665-676, Sept. 1967
2. Gaarder, N. T., "The Design of Point Detector Arrays, I," IEEE Trans. on Information Theory, vol. IT-13, No. 1, pp. 42-50, Jan. 1967.
3. Gaarder, N. T., "The Design of Point Detector Arrays, II" IEEE Trans. on Information Theory, vol. IT-12, No. 2, April 1966
4. Shor, S., "Adaptive Technique to Discriminate Against Coherent Noise in a Narrow-Band System," J. Acoustical Soc. of Am., vol. 39, No. 1, pp. 74-78, Jan. 1966
5. Mermoz, H., "Matched Filters and Optimum Use of an Array," Translation of the Proceedings-NATO Advanced Study Institute, Grenoble, France, Sept. 1964
6. Heaps, H. S., "General Theory for the Synthesis of Hydrophone Arrays," J. Acoustical Soc. Am., vol. 32, No. 3, pp. 356-363, March 1960
7. Picinbôno, B., "Optimum Filtering and Multichannel Receivers," IEEE Trans. on Information Theory, Vol. IT-12, No. 2, pp. 256-260, April 1966
8. Edelblute, D., Fisk, J., and Kinnison, C., "Criteria for Optimum Signal Detection Theory for Arrays," J. Acoustical Soc. Am., vol. 41, No. 1, pp. 199-205, Jan. 1967
9. Bryn, F., "Optimum Signal Processing of Three Dimensional Arrays Operating on Gaussian Signals and Noise," J. Acoustical Soc. Am., vol. 32, No. 3, pp. 289-297, March 1962
10. Vanderkulk, W., "Optimum Processing for Acoustic Arrays," J. Brit. IRE, vol. 26, No. 10, pp. 285-292, Oct. 1963
11. Capon, J., Greenfield, R., and Kolker R., "Multidimensional Maximum Likelihood Processing of a Large Aperture Seismic Array," Proc. IEEE, vol. 55, No. 2, pp. 192-211, Feb. 1967
12. Widrow, B., "Adaptive Filters I: Fundamentals," Stanford Electronics Laboratories Technical Report No. 6764-6, December 1966
13. Widrow, B., Mantey, P. E., Griffiths, L. J., and Goode B. B., "Adaptive Antenna Systems," Proc. IEEE, vol. 55, pp. 2143-2159, December 1967

14. Somin, M. E., "On Some Aspects of Array Matched Filtering," Ph.D. Dissertation, Polytechnic Institute of Brooklyn, June 1969
15. Griffiths, L. J., "A Simple Adaptive Algorithm for Real-Time Processing in Antenna Arrays," Proc. IEEE, vol. 57, pp. 1696-1704, October 1969
16. Lucky, R. W., "Automatic Equalization for Digital Communication," Bell System Technical J., vol. 44, No. 4, pp. 547-588, April 1965
17. Lucky, R. W., "Techniques for the Adaptive Equalization of Digital Communications Systems," Bell System Technical J., vol. 45, No. 2, pp. 255-286, Feb. 1966
18. Gersho, A., "Adaptive Equalization of Highly Dispersive Channels for Data Transmission: I," Bell System Technical J., January 1969
19. Lo, Y. T., Lee, S. W., and Lee, Q. H., "Optimization of Directivity and Signal-to-Noise Ratio of an Arbitrary Antenna Array," Proc. IEEE, vol. 54, pp. 1033-1045, August 1966
20. Gilbert, E. N., and Morgan, S. P., "Optimum Design of Directive Antenna Arrays Subject to Random Variations," Bell System Tech. J., vol. 34, pp. 637-663, May, 1955
21. Rosen, J. B., "The Gradient Projection Method for Nonlinear Programming, Part I. Linear Constraints," J.S.I.A.M., vol 8, pp. 181-217, March 1960
22. Rosen, J. B., "The Gradient Projection Method for Nonlinear Programming, Part II. Nonlinear Constraints," J.S.I.A.M., vol 9, pp. 514-532, Dec. 1961
23. Halmos, P. R. "Finite-Dimensional Vector Spaces," Second Edition, D. Van Nostrand Co., Inc. 1958
24. Zoutendijk, G., "Methods of Feasible Directions" Amsterdam: Elsevier, 1960
25. Kunzi, Krelle, Oettli, "Nonlinear Programming", Blaisdell Publishing Co., 1966
26. Hadley, G., "Nonlinear and Dynamic Programming", Addison - Wesley Publishing Co., 1964
27. Kuhn, H. W., and Tucker, A. W. "Nonlinear Programming," in "Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability," J. Neyman, editor. University of California Press, 1951, pp 481-492
28. Goldstein, A. A., "Constructive Real Analysis," Harper and Row, 1967.